

# Biased-Belief Equilibrium\*

Yuval Heller<sup>†</sup>Eyal Winter<sup>‡</sup>

February 17, 2019

## Abstract

We investigate how distorted, yet structured, beliefs can persist in strategic situations. Specifically, we study two-player games in which each player is endowed with a biased-belief function that represents the discrepancy between a player's beliefs about the opponent's strategy and the actual strategy. Our equilibrium condition requires that: (1) each player choose a best-response strategy to his distorted belief about the opponent's strategy, and (2) the distortion functions form best responses to one another, in the sense that if one of the players is endowed with a different distortion function, then that player is outperformed in the game induced by this new distortion function. We impose a mild monotonicity restriction on the feasible biased beliefs, and we obtain sharp predictions and novel insights into the set of stable outcomes and their supporting stable biases in various classes of games.

JEL classification: C73, D83.

Keywords: commitment, indirect evolutionary approach, distortions, wishful thinking, strategic complements, strategic substitutes.

## 1 Introduction

Standard models of equilibrium behavior attribute rationality to players at two different levels: beliefs and actions (see, e.g., [Aumann and Brandenburger, 1995](#)). Players are assumed to behave as if they form correct beliefs about the opponents' behavior, and they choose actions that maximize their utility given the beliefs that they hold. Much of the literature in behavioral and experimental economics that documents violations of the assumption that players have correct beliefs ascribes these violations to cognitive limitations. However, in interactive environments where one person's beliefs affect other persons' actions, belief distortions are not arbitrary, and they may arise to serve some strategic purpose.

In this paper we investigate how distorted, yet structured, beliefs can persist in strategic situations. Our basic assumption here is that distorted beliefs can persist because they offer a strategic

---

\*The authors are very grateful to the anonymous referees for very helpful comments and suggestions.

<sup>†</sup>Department of Economics, Bar Ilan University, Israel. [yuval.heller@biu.ac.il](mailto:yuval.heller@biu.ac.il). URL: <https://sites.google.com/site/yuval26/>. The author is grateful to the *European Research Council* for its financial support (ERC starting grant #677057).

<sup>‡</sup>Center for the Study of Rationality and Department of Economics, Hebrew University of Jerusalem, Israel. [mseyal@mscc.huji.ac.il](mailto:mseyal@mscc.huji.ac.il). URL: <http://www.ma.huji.ac.il/~mseyal/>. The author is grateful to the German-Israeli Foundation for Scientific Research and Google for their financial support.

advantage to those who hold them even when these beliefs are wrong. More specifically, players often hold distorted beliefs as a form of commitment device that affects the behavior of their counterparts. The precise cognitive process that is responsible for the formation of beliefs is complex, and it is beyond the scope of this paper to outline it. We believe, however, that in addition to analytic assessment of evidence, preferences in the form of desires, fears, and other emotions contribute to the process and, to an extent, facilitate belief biases. If the evidence is unambiguous and decisive, or if the consequence of belief distortion is detrimental to the player’s welfare, preferences may play less of a role and learning may work to calibrate beliefs to reality. But when beliefs are biased in ways that favor their holders by affecting the behavior of their counterparts, learning can actually reinforce biases rather than diminish them.

**Biased Beliefs** Standard equilibrium notions in game theory draw a clear line between preferences and beliefs. The former are exogenous and fixed; the latter can be amended through Bayesian updating but are not allowed to be affected by preferences. However, phenomena such as wishful thinking (see, e.g., [Babad and Katz, 1991](#)) and overconfidence (see, e.g., [Forbes, 2005](#); [Barber and Odean, 2001](#); [Malmendier and Tate, 2005](#); [Heller, 2014](#)), where beliefs are tilted toward what their holder desires reality to be, suggest that in real life, beliefs and preferences can intermingle, and that biased beliefs may be persistent. Similarly, belief rigidity and belief polarization (see, e.g., [Lord, Ross, and Lepper, 1979](#); [Ross and Anderson, 1982](#)) refer to situations in which two people with conflicting prior beliefs each strengthen their beliefs in response to observing the same data. The parties’ aversion to depart from their original beliefs can also be regarded as a form of interaction between preferences and beliefs.

It is easy to see how the belief biases described above can have strategic benefits in interactive situations. Wishful thinking and optimism can facilitate cooperation in interactions that require mutual trust. Overconfidence can deter competitors, and belief rigidity can allow an agent to support a credible threat. An important objective of our analysis is to identify the strategic environments that support biases such as wishful thinking as part of equilibrium behavior. It is worthwhile to note that individuals are not the only ones susceptible to strategically motivated belief biases. Governments are prone to be affected by such biases as well. The Bush administration’s unsubstantiated confidence in Saddam Hussein’s possession of “weapons of mass destruction” prior to the Second Gulf War and the vast discrepancy between Israeli and US intelligence assessments of Iran’s nuclear intentions prior to the signing of the Iran nuclear deal can be interpreted as strategically motivated belief distortion.<sup>1</sup>

Belief biases in strategic environments are also connected to self-interest biases regarding moral and ethical standards. [Babcock and Loewenstein \(1997\)](#) had participants in a lab experiment ne-

---

<sup>1</sup>There are other possible interpretations of these controversial real-life examples. In a dynamic real-life setup it is hard to have access to agents’ private information, and therefore it is very difficult to achieve direct empirical evidence for persistent biased beliefs. There are a few lab experiments that elicit subjects’ beliefs (using monetary incentives and proper scoring rules) about the expected behavior of the opponent. [Nyarko and Schotter \(2002\)](#) demonstrate that the elicited forecasts of subjects about the opponents’ future behavior substantially differ from the empirical play of opponents in the past. [Palfrey and Wang \(2009\)](#) present evidence that forecasts by players (about the opponent’s behavior in a simple two-player game) are significantly different from the forecasts of external observers. Moreover, the players’ forecasts are systematically biased, and significantly less accurate than the forecasts of the external observers.

gotiate a deal between a plaintiff and a defendant in a court case. When they asked participants to make predictions about the outcome of the real court case the authors found a significant belief divergence depending on the role participants were assigned to in the negotiations. A similar moral hypocrisy was revealed by [Rustichini and Villeval \(2014\)](#) who showed that subjects’ subjective judgments regarding fairness in bargaining depended on the bargaining power they were assigned in the experiment.

A different body of empirical evidence consistent with strategic beliefs is offered by the psychiatric literature on “depressive realism” (e.g., [Dobson and Franche, 1989](#)). This literature compares probabilistic assessments conveyed by psychiatrically healthy people with those suffering from clinical depression. Participants in both categories were requested to assess the likelihood of experiencing negative or positive events in both public and private setups. Comparing subjects’ answers with the objective probabilities of these events revealed that in a public setup clinically depressed individuals were more realistic than their healthy counterparts for both types of events. The apparent belief bias among healthy individuals can be reasonably attributed to the strategic component of beliefs. Mood disorders negatively affect strategic reasoning ([Inoue, Tonooka, Yamada, and Kanba, 2004](#)), which, to a certain extent, may diminish strategic belief distortion among clinically depressed individuals relative to their healthy counterparts.

For biased beliefs to yield a strategic advantage to the agents holding them, it is essential that (1) agents be committed to follow their biased beliefs, and (2) agents best-reply to the perceived behavior induced by their counterparts’ biases (both on and off the equilibrium path). For the sake of tractability, we shall avoid formalizing a concrete dynamic model that describes how biased beliefs are formed, and how agents credibly commit to these biased beliefs. Instead, we shall adopt a static approach by imposing equilibrium conditions on the agents’ beliefs and the opponents’ interpretation of their beliefs. (We discuss our modeling approach and its evolutionary interpretation in [Section 3.6](#), and we present a formal evolutionary foundation in [Appendix B](#).) This static approach is consistent with a large part of the literature on endogenous preferences (see, e.g., the literature cited below). Nevertheless, we mention a few mechanisms that can facilitate these processes and turn biased beliefs into a credible commitment device.

1. Refraining from accessing or using biased sources of information, e.g., subscribing to a newspaper with a specific political orientation, consulting biased experts, and reading Facebook’s personalized news feeds, which are typically biased due to friends who hold similar beliefs.
2. Passionately following a religion, a moral principle, or an ideology that has belief implications on human behavior.
3. Possessing personality traits that have implications on beliefs (e.g., narcissism or naivety).

The mechanisms described above are likely not only to induce belief biases, but also to generate signals sent to the player’s counterparts about these biases with a certain degree of verifiability. These mechanisms, the signals they induce, and their interpretation are the main forces that facilitate biased-belief equilibrium.

**Solution Concept** Our notion of biased-belief equilibrium (henceforth, BBE) uses a two-stage paradigm. In the first stage each player is endowed with a biased-belief function. This function represents the discrepancy between a player’s beliefs about the strategy profile of other players and the actual profile. In the second stage the players play the biased game induced by their distortion functions, in which each player chooses a best-reply strategy to his biased belief about the opponent’s strategy (the chosen strategy profile is referred to as the equilibrium outcome). Finally, our equilibrium condition requires that the distortion functions not be arbitrary, but form best replies to one another.

If one of the players deviates to being endowed with a different biased-belief function, then there might be multiple Nash equilibria in the new biased game induced by this deviation. Our weak notion (*weak BBE*) requires the deviator to be outperformed in at least one equilibrium of the new biased game. Our strong notion (*strong BBE*) requires (1) each agent to have a *monotone* biased belief, according to which he assigns a higher probability to his opponent playing a certain strategy than this probability actually is, and (2) a deviator to be outperformed in *all* Nash equilibria of the new biased game. Our main notion, *BBE*, lies in between these two notions, and it requires (1) each agent to have a monotone biased belief, and (2) the deviator to be outperformed in at least one plausible Nash equilibrium of the new biased game, where we rule out implausible Nash equilibria in which the non-deviator behaves differently even though he does not observe any change in the deviator’s perceived strategy.

In Section 2.5 we present our main evolutionary interpretation of our solution concept, according to which the endowed biased beliefs are the result of an evolutionary process of social learning (the interpretation is formalized in Appendix B). In addition, we present an alternative, delegation interpretation of the model (which is formalized in Appendix C).

**Nash Equilibrium and BBE** We begin our analysis by studying the relations between BBE outcomes and Nash equilibria. We show that any Nash equilibrium can be implemented as the outcome of a BBE, though in some cases this requires that the players have biased beliefs that are accurate on the equilibrium path, but that they be blind to some deviations of the opponent off the equilibrium path. This, in particular, implies that every game admits a BBE. Next, we show that introducing biased beliefs does not change the set of equilibrium outcomes in games in which at least one of the players has a dominant action. By contrast, BBE admits non-Nash behavior in most other games.

**Main Results** Our main results show that the notion of BBE induces substantial predictive power in various classes of interval games. In these classes of games the strategy of each player is a number in a bounded interval, where a higher strategy (interpreted as a higher investment) induces a higher payoff for the opponent. We begin by characterizing the set of BBE in games with strategic complements (Bulow, Geanakoplos, and Klemperer, 1985), such as price competition with differentiated goods (Example 2), input games (Example 9 in Appendix A.4), and stag hunt games (Example 10 in Appendix A.4). We show three key properties of any BBE: (1) *overinvestment*: the strategy of each agent is (weakly) higher than the best reply to the opponent’s (real) strategy,

(2) *ruling out bad outcomes*: both players invest more than their investments in the worst Nash equilibrium of the underlying game, and (3) *wishful thinking*: each agent perceives his opponent as investing (weakly) more than the opponent’s real investment.

Next, we characterize the set of BBE in games with strategic substitutes, such as Cournot competitions (Example 3) and hawk-dove games (Example 5 in Appendix A.5). We show three key properties of any BBE: (1) *underinvestment*: the strategy of each agent is (weakly) higher than the best reply to the opponent’s (real) strategy, (2) *ruling out excellent outcomes*: at least one of the players invests less than his investments in one of the Nash equilibria of the underlying game, and (3) *wishful thinking*: each agent perceives his opponent as investing (weakly) more than the opponent’s real investment.

Finally, we characterize the set of BBE in a class of games (which are less common in economic interactions), in which the strategy of player 1 is a complement of player 2’s strategy, while the strategy of player 2 is a substitute of player 1’s strategy (e.g., duopolistic competition in which one firm chooses its quantity while the opposing firm chooses its price (Singh and Vives, 1984), and various classes of asymmetric contests (Dixit, 1987)). We show that in this class of games agents present *pessimism* in any BBE: each agent perceives his opponent as investing (weakly) less than the opponent’s real investment.

**Additional Results** Our next result shows an interesting class of BBE that exist in all games. We say that a strategy is undominated Stackelberg if it maximizes a player’s payoff in a setup in which the player can commit to an undominated strategy, and his opponent reacts by best-replying to this strategy. We show that every game admits a BBE in which one of the players is “strategically stubborn” in the sense of having a constant belief about the opponent’s strategy, and always playing his undominated Stackelberg strategy, while the opponent is “rational” in the sense of having undistorted beliefs and best-replying to the player’s true strategy.

Section 7.2 shows that unless one imposes both requirements on the definition of a BBE, namely, monotonicity and ruling out implausible equilibria, then the set of BBE outcomes is very large in various classes of games. Specifically, Proposition 8 shows that for a large class of finite games, a strategy profile is a monotone weak BBE iff (1) no player uses a strictly dominated strategy, and (2) the payoff of each player is above the minmax payoff of the player in a setup in which both players are restricted to choose only undominated strategies (i.e., strategies that are not strictly dominated). Proposition 9 shows a similar folk theorem result for non-monotone strong BBE in a large class of interval games.

**Empirical Predictions** Our main results imply two empirical predictions. First, they suggest that efficient (non-Nash equilibrium) outcomes are easier to support in games with strategic complements, relative to games with strategic substitutes. This prediction is consistent with the experimental findings of Potters and Suetens (2009), which show that there is significantly more cooperation in games with strategic complements than in the case of strategic substitutes.

Our second empirical prediction is that wishful thinking is strategically stable in many common environments, though some (less common) strategic interactions may induce pessimism. This

empirical prediction is consistent with the experimental evidence that people tend to present wishful thinking, while the presented level of wishful thinking may substantially differ between various environments; see, e.g., [Babad and Katz \(1991\)](#); [Budescu and Bruderman \(1995\)](#); [Bar-Hillel and Budescu \(1995\)](#) and [Mayraz \(2013\)](#).

**Structure** The structure of this paper is as follows. We discuss the related literature in Section 2. Section 3 describes the model. In Section 4 we analyze the relations between BBE and Nash equilibria. Section 5 defines games with strategic complements/substitutes and wishful thinking. We analyze these games and present our main results in Section 6. In Section 7 we present additional results: (1) the relation between BBE and strategies played by a Stackelberg leader, and (2) folk theorem results when relaxing the definition of BBE. We conclude in Section 8. All the appendices of the paper appear in the online supplementary material. Appendix A presents various interesting examples. We formally present the evolutionary interpretation of our solution concept in Appendix B, and the delegation interpretation in Appendix C. Appendix D relaxes the assumption that biased beliefs have to be continuous. Appendix E shows how to extend our results to a setup with partial observability. Appendix F presents our formal proofs.

## 2 Related Literature and Contributions

Our paper aims at making a contribution to the behavioral game theory literature. Much of this literature concerns behavioral equilibrium concepts that depart from the framework of Nash equilibrium by introducing weaker rationality conditions. This has been done primarily at the level of preferences (e.g., [Güth and Yaari, 1992](#); [Fehr and Schmidt, 1999](#); [Bolton and Ockenfels, 2000](#); [Acemoglu and Yildiz, 2001](#); [Heifetz, Segev, et al., 2004](#); [Dekel, Ely, and Yilankaya, 2007](#); [Heifetz, Shannon, and Spiegel, 2007a](#); [Friedman and Singh, 2009](#); [Herold and Kuzmics, 2009](#); [Heller and Winter, 2016](#); [Winter, Garcia-Jurado, and Mendez-Naya, 2017](#)). But it has also been done at the level of beliefs (e.g., [Geanakoplos, Pearce, and Stacchetti, 1989](#); [Rabin, 1993](#); [Battigalli and Dufwenberg, 2007](#); [Attanasi and Nagel, 2008](#); [Battigalli and Dufwenberg, 2009](#); [Battigalli, Dufwenberg, and Smith, 2015](#); [Gannon and Zhang, 2017](#)). This latter literature deals with belief-dependent preferences, and focuses primarily on the way players' beliefs about the intentions of others affect their preferences and behavior.

Our equilibrium concept also operates on beliefs rather than preferences but is based on an inherently different approach. Preferences in our model are not affected by beliefs but beliefs are biased in a way that serves players' strategic purposes. Our analysis of biased belief goes beyond characterizing equilibrium outcomes. An additional important objective is to identify the belief biases that support these equilibrium outcomes in different strategic environments. Central to our analysis are belief-distortion properties, such as wishful thinking and pessimism, that sustain BBE in different strategic environments.

The existing literature has presented various prominent solution concepts that assume that players have distorted beliefs. Some examples include models of level- $k$  and cognitive hierarchy (see, e.g., [Stahl and Wilson, 1994](#); [Nagel, 1995](#); [Costa-Gomes, Crawford, and Broseta, 2001](#); [Camerer, Ho, and](#)



Chong, 2004), analogy-based expectation equilibrium (Jehiel, 2005), cursed equilibrium (Eyster and Rabin, 2005), and Berk-Nash equilibrium (Esponda and Pouzo, 2016). These equilibrium notions have been helpful in understanding strategic behavior in various setups, and yet these notions pose a conceptual challenge to our understanding of the persistence of distorted beliefs, even in view of the empirical evidence for such persistence. If players can infer the truth *ex post* why don't they calibrate their beliefs toward reality? Much of the literature presenting such models points to cognitive limitations as the source of this rigidity. Our model and analysis offer an additional perspective to this issue by suggesting that belief biases that yield a strategic advantage in the long run are likely to emerge in equilibrium. In this sense our approach can be viewed as providing a tool to explain why some cognitive limitations persist while others do not (see Example 11 in Appendix A, in which we show how level-1 behavior can be supported as part of a BBE outcome in the traveler's dilemma).

Our notion of BBE is related to the notion of conjectural equilibrium (Battigalli and Guaitoli, 1997, originally written in 1988) insofar as both solution concepts relax the Nash equilibrium's requirement that beliefs need to be consistent with actual play (while still requiring that an agent's action has to be optimal given the agent's belief). A conjectural equilibrium is defined in an environment in which players do not observe each other's actions but rather observe signals of each other's actions, according to an exogenous feedback correspondence. In a conjectural equilibrium each player best replies to his belief about the opponent's action, and this belief is required to be consistent with the signal observed by the player. There are two key structural differences between a BBE and a conjectural equilibrium. First, a BBE is defined in an environment in which there is no exogenous feedback correspondence; rather, the feedback correspondence is implicitly defined as part of the solution concept by the agents' biased-belief functions. These biased-belief functions are not restricted by a consistency requirement with respect to an exogenous feedback mechanism, but rather they are restricted by the requirement that each biased-belief function has to be a best reply against the opponent's biased belief. The second structural difference is that while a BBE describes what would be the agent's belief for *any* feasible action of the opponent, a conjectural equilibrium describes only the agent's belief about the equilibrium action of the opponent.

Despite these structural differences, it is interesting to discuss relations between the equilibrium behavior induced by each solution concept, i.e., the relations between a BBE outcome and a conjectural equilibrium outcome. Without restricting the feedback correspondence, the notion of conjectural equilibrium is rather broad (it rules out only strictly dominated strategies), and, accordingly, any BBE outcome is a conjectural equilibrium outcome. Fudenberg and Levine's (1993) notion of self-confirming equilibrium deals with extensive-form games, and refines conjectural equilibrium by requiring that the feedback correspondence is the one in which each player observes the opponent's realized actions (but does not observe the opponent's behavior off the equilibrium path). In the setup of two-player one-shot games, which is the focus of the present paper, the set of self-confirming equilibria coincides with the set of Nash equilibria (whereas the set of BBE outcomes is broader and includes non-Nash outcomes). Another refinement of conjectural equilibrium is the rationalizable conjectural equilibrium (Rubinstein and Wolinsky, 1994; the notion has been generalized to games with structural uncertainty in Esponda, 2013). This concept requires that the

agents' beliefs be consistent with the common knowledge that all players maximize utility given their signals. There is no inclusion relation between the set of BBE outcomes and the set of rationalizable conjectural equilibrium outcomes. Specifically, in games with a unique rationalizable action profile, such as price competitions with differentiated goods and Cournot competitions, the unique rationalizable conjectural equilibrium outcome is the Nash equilibrium (for any feedback correspondence), while the set of BBE outcomes is substantially larger (see Examples 2 and 3). By contrast, in games such as stag hunt and hawk–dove, when the feedback correspondence is non-informative any action profile is a conjectural equilibrium outcome, while the set of BBE outcomes is much more restricted (see Examples 10 and 12 in Appendix A).

### 3 Model

#### 3.1 Underlying Game

Let  $i \in \{1, 2\}$  be an index used to refer to one of the players in a two-player game, and let  $j$  be an index referring to the opponent. Let  $G = (S, \pi)$  be a normal-form two-player game (henceforth, *game*), where  $S = (S_1, S_2)$  and each  $S_i$  is a convex compact set of strategies. Specifically, we focus on two cases:

1. *Finite games*: Each  $S_i$  is a simplex over a finite set of pure actions, where each strategy corresponds to a mixed action (i.e.,  $A_i$  is a finite set of pure actions, and  $S_i = \Delta(A_i)$ ), and the von Neumann–Morgenstern payoff function is linear with respect to the mixing probability.
2. *Interval games*: Each  $S_i$  is a bounded interval in  $\mathbb{R}$  (e.g., each player chooses a real number representing quantity, price, or effort).

We denote by  $\pi = (\pi_1, \pi_2)$  players' payoff functions; i.e.,  $\pi_i : S \rightarrow \mathbb{R}$  is a function assigning each player a payoff for each strategy profile. We use  $s_i$  to refer to a typical strategy of player  $i$ . We assume each payoff function  $\pi_i(s_i, s_j)$  to be continuously twice differentiable in both parameters and weakly concave in the first parameter ( $s_i$ ).

Let  $BR$  (resp.,  $BR^{-1}$ ) denote the (inverse) best-reply correspondence; i.e.,

$$BR(s_i) = \operatorname{argmax}_{s_j \in S_j} (\pi_j(s_i, s_j))$$

is the set of best replies against strategy  $s_i \in S_i$ , and

$$BR^{-1}(s_i) = \{s_j \in S_j | s_i \in BR(s_j)\}$$

is the set of strategies for which  $s_i$  is a best reply against them.

In a finite game, we use  $a_i \in A_i$  to denote also the degenerate mixed action that assigns mass one to  $a_i$ . When the set of actions of a player is given as an ordered set  $A_i = (a_i^1, a_i^2, \dots, a_i^n)$ , we identify a mixed action with a vector  $s_i = (\alpha_1, \alpha_2, \dots, \alpha_n)$ , where  $0 \leq \alpha_k = s_i(a_i^k)$  for each  $1 \leq k \leq n$ , and  $\sum_k \alpha_k = 1$ . Given two strategies  $s_i, s'_i \in S_j$  and  $\alpha \in [0, 1]$ , let  $\alpha \cdot s_i + (1 - \alpha) \cdot s'_i$  be the mixture of the two strategies:  $(\alpha \cdot s_i + (1 - \alpha) \cdot s'_i)(a_i) = \alpha \cdot s_i(a_i) + (1 - \alpha) \cdot s'_i(a_i)$ .



When there are two (ordered) actions for each player (say,  $A_i = \{c_i, d_i\}$ ), we identify a mixed action  $s_i$  with the probability it assigns to the first pure action  $s_i(c_i)$ , and we identify the set of strategies  $S_i$  with the interval  $[0, 1]$ . Thus, a game with two actions for each player can be captured both as a finite game and as an interval game.

### 3.2 Biased-Belief Function

We start here with the definition of biased-belief functions that describe how players' beliefs are distorted. A *biased belief*  $\psi_i : S_j \rightarrow S_j$  is a *continuous* function that assigns to each strategy of the opponent, a (possibly distorted) belief about the opponent's play. That is, if the opponent plays  $s_j$ , then player  $i$  believes that the opponent plays  $\psi_i(s_j)$ . We call  $s_j$  the opponent's real strategy, and we call  $\psi_i(s_j)$  the opponent's perceived (or biased) strategy. Formally, the continuity requirement is that if  $(s_{j,n})_n \xrightarrow{n \rightarrow \infty} s_j$ , then  $(\psi_i(s_{j,n}))_n \xrightarrow{n \rightarrow \infty} \psi_i(s_j)$  (where in a finite game, we say that  $(s_{j,n})_n \xrightarrow{n \rightarrow \infty} s_j$  iff  $(s_{j,n}(a))_n \xrightarrow{n \rightarrow \infty} s_j(a)$  for each action  $a$ ).

*Remark 1.* Two reasons motivate us to require that a biased belief be continuous: (1) continuity implies that each biased game (defined below) admits a Nash equilibrium, which allows us to simplify the definition of BBE, and (2) continuity reflects a plausible restriction that a small change in the opponent's strategy should induce a small change in the perceived strategy. In Appendix D we present an alternative (and somewhat more complicated) definition of a BBE that relaxes the assumption that biased beliefs must be continuous, and we show that all the BBE characterized in the results of the paper remain BBE when we allow deviators to use discontinuous biased beliefs.

We say that a biased belief  $\psi_i : S_j \rightarrow S_j$  is *monotone* if:

1. In interval games:  $s_j \geq s'_j$  implies  $\psi_i(s_j) \geq \psi_i(s'_j)$  for each strategy  $s_j \in S_j$ .
2. In finite games: If the opponent plays  $a_j$  more often, while keeping the same proportion of playing the remaining actions, then the perceived probability that the opponent plays any other action weakly decreases (which implies, in particular, that the perceived probability that the opponent plays  $a_j$  weakly increases); that is,

$$(\psi_i((1 - \alpha) \cdot s_j + \alpha \cdot a_j)) (a'_j) \leq (\psi_i(s_j)) (a'_j)$$

for each  $\alpha \in [0, 1]$ , each action  $a_j \in A_j$ , each action  $a'_j \neq a_j$ , and each strategy  $s_j \in \Delta(A_j)$ . In particular, when the game has two actions for each player, a biased belief  $\psi_i$  is monotone iff  $\psi_i$  is weakly increasing in  $\alpha_j$ ; i.e.,  $\alpha_j \geq \alpha'_j$  implies that  $\psi_i(\alpha_j) \geq \psi_i(\alpha'_j)$ .

Monotone biased beliefs reflect a plausible restriction on the distortion of agents, namely, that if the opponent changes his real strategy in some direction, the agent captures the direction of the change correctly, but may have the wrong perception about the magnitude of the change.

Let  $I_d$  be the undistorted (identity) function, i.e.,  $I_d(s) = s$  for each strategy  $s$ . A biased belief  $\psi$  is *blind* if the perceived opponent's strategy is independent of the opponent's real strategy, i.e., if  $\psi(s_j) = \psi(s'_j)$  for each  $s_j, s'_j \in S_j$ . With a slight abuse of notation we use  $s_i$  to denote also the blind biased belief  $\psi_j$  that is always equal to  $s_i$ .

### 3.3 Biased Game

An underlying game and a profile of biased beliefs jointly induce a biased game in which the (biased) payoff of each player is determined by the perceived strategy of the opponent. Formally:

**Definition 1.** Given an underlying game  $G = (S, \pi)$  and a profile of biased beliefs  $(\psi_i, \psi_j)$ , let the *biased game*  $G_\psi = (S, \psi \circ \pi)$  be defined as the game with the following payoff function  $(\psi \circ \pi)_i : S_i \times S_j \rightarrow \mathbb{R}$  for each player  $i$ :

$$(\psi \circ \pi)_i(s_i, s_j) = \pi_i(s_i, \psi_i(s_j)).$$

A Nash equilibrium of a biased game is defined in the standard way. Formally, a pair of strategies  $s^* = (s_1^*, s_2^*)$  is a Nash equilibrium of a biased game  $G_\psi = (S, \psi \circ \pi)$ , if each  $s_i^*$  is a best reply against the perceived strategy of the opponent, i.e.,

$$s_i^* = \operatorname{argmax}_{s_i \in S_i} \left( \pi_i \left( s_i, \psi_i(s_j^*) \right) \right).$$

Let  $NE(G_\psi) \subseteq S_1 \times S_2$  denote the set of all Nash equilibria of the biased game  $G_\psi$ .

Observe that the set of strategies of a biased game is convex and compact, and the payoff function  $(\psi \circ \pi)_i : S_i \times S_j \rightarrow \mathbb{R}$  is weakly concave in the first parameter and continuous in both parameters. This implies (due to a standard application of Kakutani's fixed-point theorem) that each biased game  $G_\psi$  admits a Nash equilibrium (i.e.,  $NE(G_{(\psi_i', \psi_j^*)}) \neq \emptyset$ ).

### 3.4 Weak and Strong BBE

We are now ready to define our equilibrium concept. A weak biased-belief equilibrium (abbr. weak BBE) is a pair consisting of a profile of biased beliefs and a profile of strategies, such that: (1) each strategy is a best reply to the perceived strategy of the opponent, and (2) each biased belief is a best reply to the opponent's biased belief, in the sense that any agent who chooses a different biased-belief function is outperformed in at least one equilibrium in the new biased game (relative to the agent's payoff in the original equilibrium). Formally:

**Definition 2.** A *weak BBE* is a pair  $(\psi^*, s^*)$ , where  $\psi^* = (\psi_1^*, \psi_2^*)$  is a profile of biased beliefs and  $s^* = (s_1^*, s_2^*)$  is a profile of strategies satisfying: (1)  $(s_i^*, s_j^*) \in NE(G_{\psi^*})$ , and (2) for each player  $i$  and each biased belief  $\psi_i'$ , there exists a strategy profile  $(s_i', s_j') \in NE(G_{(\psi_i', \psi_j^*)})$ , such that the following inequality holds:  $\pi_i(s_i', s_j') \leq \pi_i(s_i^*, s_j^*)$ .

The notion of weak BBE is arguably too permissive because it allows incumbents: (1) to have implausible non-monotone beliefs, and (2) to outperform the deviators in a single Nash equilibrium of the biased game (while, possibly, the incumbents are outperformed by the deviators in many other equilibria). Proposition 8 (in Section 7.2.2) demonstrates that this single Nash equilibrium, in which the deviators are outperformed, may be implausible due to allowing the incumbents to “discriminate” against the deviators, even though the deviators exhibit exactly the same perceived behavior as the rest of the population.

The more restrictive refinement of strong BBE requires that (1) incumbents have monotone beliefs, and (2) deviators who choose a different biased-belief function be outperformed in *all* equilibria of the induced biased game. Formally:

**Definition 3.** A *weak BBE*  $(\psi^*, s^*)$  is a *strong BBE* if (1) each biased function  $\psi_i^*$  is monotone, and (2) the inequality  $\pi_i(s'_i, s'_j) \leq \pi_i(s_i^*, s_j^*)$  holds for every player  $i$ , every biased belief  $\psi'_i$ , and every strategy profile  $(s'_i, s'_j) \in NE(G_{(\psi'_i, \psi_j^*)})$ .

### 3.5 BBE

Finite games typically induce multiple Nash equilibria. This is often the case also with respect to biased games. This suggests that the refinement of strong BBE may be too restrictive, as there are potentially many Nash equilibria of many biased games, and the requirement of the deviators being outperformed in all these equilibria might be too demanding. Our main solution concept, BBE, lies in between weak BBE and strong BBE.

In a BBE, the deviator is required to be outperformed in at least one *plausible* equilibrium of the new biased game. Roughly speaking, in a plausible equilibrium of the new biased game induced by a deviation of player  $i$  to a different biased belief, player  $j$  is allowed to choose a new strategy only if he distinguishes between  $i$ 's original strategy and  $i$ 's new strategy. More precisely, implausible equilibria are defined as follows. We say that a Nash equilibrium of a biased game induced by a deviation of player  $i$  is implausible if (1) player  $i$ 's strategy is perceived by the non-deviating player  $j$  as coinciding with player  $i$ 's original strategy, (2) player  $j$  plays differently relative to his original strategy, and (3) player  $j$  playing his original strategy induces an equilibrium of the biased game. That is, implausible equilibria are those in which the non-deviating player  $j$  plays differently against a deviator even though player  $j$  has no reason to do so: player  $j$  does not observe any change in player  $i$ 's behavior, and player  $j$ 's original behavior remains an equilibrium of the biased game. Formally:

**Definition 4.** Given weak BBE  $(\psi^*, s^*)$ , deviating player  $i$ , and biased belief  $\psi'_i$ , we say that a Nash equilibrium of the biased game  $(s'_i, s'_j) \in NE(G_{(\psi'_i, \psi_j^*)})$  is *implausible* if: (1)  $\psi_j^*(s'_i) = \psi_j^*(s_i^*)$ , (2)  $s_j^* \neq s'_j$ , and (3)  $(s'_i, s_j^*) \in NE(G_{(\psi'_i, \psi_j^*)})$ . An equilibrium is *plausible* if it is not implausible. Let  $PNE(G_{(\psi'_i, \psi_j^*)})$  be the set of all plausible equilibria of the biased game  $G_{(\psi'_i, \psi_j^*)}$ .

Note that it is immediate from Definition 4 and the nonemptiness of  $NE(G_{(\psi'_i, \psi_j^*)})$  that  $PNE(G_{(\psi'_i, \psi_j^*)})$  is nonempty.

**Definition 5.** Weak BBE  $(\psi^*, s^*)$  is a *BBE* if (1) each biased function  $\psi_i^*$  is monotone, and (2) for each player  $i$  and each biased belief  $\psi'_i$ , there exists a plausible Nash equilibrium  $(s'_i, s'_j) \in PNE(G_{(\psi'_i, \psi_j^*)})$ , such that  $\pi_i(s'_i, s'_j) \leq \pi_i(s_i^*, s_j^*)$ .

A strategy profile  $s^* = (s_1^*, s_2^*)$  is a (*resp.*, *strong*, *weak*) *BBE outcome* if there exists a profile of biased beliefs  $\psi^* = (\psi_1^*, \psi_2^*)$  such that  $(\psi^*, s^*)$  is a (*resp.*, *strong*, *weak*) BBE. In this case we say that the biased belief  $\psi^*$  supports (or implements) the outcome  $s^*$ .

### 3.6 Discussion of the Model

**Evolutionary/Learning Interpretation** Biases can emerge in a learning process that reinforces biases that yield a strategic advantage to their holders. Specifically, we interpret a BBE to be a reduced-form solution concept capturing the essential features of an evolutionary process of cultural or social learning. Our methodology follows the extensive literature that studies the stability of endogenous preferences using the “indirect evolutionary approach” (see, e.g., Güth and Yaari, 1992; Güth, 1995; Fershtman and Weiss, 1998; Dufwenberg and Güth, 1999; Koçkesen, Ok, and Sethi, 2000; Guttman, 2003; Güth and Napel, 2006; Heifetz, Shannon, and Spiegel, 2007b; Friedman and Singh, 2009; Herold and Kuzmics, 2009; Alger and Weibull, 2013; Heller and Mohlin, 2017). We apply this modeling approach to the study of endogenous biased beliefs in a setup in which biased beliefs induce behavior, behavior determines “success,” and success regulates the evolution of biased beliefs.

In Appendix B we formally adapt the definition of a stable population state from Dekel, Ely, and Yilankaya (2007) to our setup, and show that the adapted definition is equivalent to a strong BBE. In what follows we briefly and informally present our evolutionary interpretation. Consider two large populations of agents: agents who play the role of player 1, and agents who play the role of player 2. In each round agents from each population are randomly matched to play a two-person game against opponents from the other population. Each agent in each population is endowed with a biased-belief function. For simplicity, we focus on “homogeneous” populations, in which all agents in the population have the same monotone biased-belief function. Agents distort their perception about the behavior of the agents in the other population according to their endowed biased-belief functions, and they play a Nash equilibrium of the biased game.

With small probability a few agents (“mutants”) in one of the populations (say, population 1) may be endowed with a different biased-belief function due to a random error or experimentation. We assume that agents of population 2 observe whether their opponents are mutants or not, and that the agents of population 2 and the mutants of population 1 gradually adapt their play against each other into an equilibrium of the new biased game. Note that a dynamic adaptation into playing a Nash equilibrium of the biased game requires agents of population 2 to know the perceived strategy currently being played by the mutants of population 1, but the agents do not need to know the biased beliefs of the mutants of population 1.

Finally, we assume that the total “success” (fitness) of agents is monotonically influenced by their (unbiased) payoff in the underlying game, and that there is a slow process in which the composition of the population evolves. This slow process might be the result of a slow flow of new agents who join the population. Each new agent randomly chooses one of the incumbents in his own population as a “mentor” (and mimics the mentor’s biased belief), where the probabilities are such that agents with higher fitness are more likely to be chosen as mentors. If the original population state is not a BBE, it implies that there are mutants who outperform the remaining incumbents in their own

population, which in turn implies that the original population state is not stable, as new agents are likely to mimic more successful mutants. By contrast, if the original population state is a BBE, it implies that for any mutant there is a new equilibrium in which the mutants are weakly outperformed relative to the incumbents of their own population, and this can allow the BBE to remain a stable state (as illustrated in the detailed example in Appendix B.3).

### Variants of the Solution Concept

*The main solution concept we use in the paper is BBE.*

In Section 7.2 we demonstrate that unless one applies both requirements of Definition 5, namely, monotonicity and ruling out implausible equilibria, then the set of BBE is very large (folk theorem results), and some of the biased beliefs that support some of these equilibria seem implausible. The intuition for the monotonicity requirement is quite straightforward (ruling out peculiar biased beliefs in which an opponent who deviates to play a higher strategy is perceived as deviating to play a lower strategy). The second requirement rules out implausible equilibria in which a player responds to his opponent's deviation in spite of not being able to perceive it

In what follows we sketch a dynamic justification for the second requirement of ruling out implausible equilibria (following the evolutionary interpretation described above). Consider a BBE  $((\psi_1^*, \psi_2^*), (s_1^*, s_2^*))$ . Assume that both  $(s'_1, s'_2)$  and  $(s'_1, s_2^*)$  are Nash equilibria of the biased game  $G_{(\psi_1', \psi_2^*)}$ . In what follows, we briefly, and informally, explain why  $(s'_1, s'_2)$  is not a plausible equilibrium of the new biased game (and, thus, why it is ruled out in the definition of BBE). Consider a deviation of some agents in the population playing in the role of player 1 to having the biased belief  $\psi_1'$ . Following this deviation, strategy  $s_1^*$  might not be a best reply to the perceived strategy of player 2 (i.e.,  $s_1^* \notin BR(\psi_1'(s_2^*))$ ) and, as a result, the deviating agents might change their strategy to  $s'_1$ , which is a best reply to the perceived strategy of player 2 (i.e.,  $s'_1 \in BR(\psi_1'(s_2^*))$ ). The current strategy profile  $(s'_1, s_2^*)$  is a Nash equilibrium of the biased game (i.e.,  $(s'_1, s_2^*) \in NE(G_{(\psi_1', \psi_2^*)})$ ). In order to move from this equilibrium to  $(s'_1, s'_2)$ , agents of population 2, who are matched against the deviators, have to change their behavior from  $s_2^*$  to  $s'_2$ , but there is no reason for them to do so, as their current behavior (namely,  $s_2^*$ ) is already a best reply to the perceived strategy of the deviators (i.e.,  $s_2^* \in BR(\psi_1^*(s'_1))$ ), as well as being how they are used to playing against non-deviators.

### Delegation Interpretation

A different interpretation of our solution concept relies on strategic delegation. The literature on strategic delegation (see, e.g., Fershtman, Judd, and Kalai, 1991; Dufwenberg and Güth, 1999; Fershtman and Gneezy, 2001) deals with players who strategically use other agents to play on their behalf, where the agents so used may have different preferences than the players using them. We adapt this approach to our setup in which agents differ in their biased beliefs (rather than in their preferences). Specifically, in Appendix C we show that the notion of weak BBE is equivalent to a subgame-perfect equilibrium of a two-stage game in which in stage one each unbiased player strategically chooses the biased belief of his agent, and in the second stage the biased agents play on behalf of the players (and each agent can observe the opposing agent's biased beliefs).

**Partial Observability** The requirement that an agent be able to observe that his opponent belongs to a group of “mutant” agents who have different biased beliefs than the rest of the population can be explained by pre-play social cues and messages that facilitate this observation. In Appendix E we show that this observability need not be perfect. We generalize the model to partial observability by studying a setup in which, when an agent is matched with a mutant opponent, the agent privately observes the opponent to be a mutant with probability  $0 < p \leq 1$ . We show that all our results hold in this extended setup for  $p$  sufficiently close to one (and some of the results hold also for low levels of  $p$ ).

## 4 Nash Equilibria and BBE Outcomes

In this section we study the relations between Nash equilibria and BBE outcomes.

### 4.1 Nash Equilibria and Biased Beliefs

We begin with a simple observation that shows that in any weak BBE in which the outcome is not a Nash equilibrium, at least one of the players must distort the opponent’s perceived strategy. The reason for this observation is that if both players have undistorted beliefs, then it must be that each agent best-responds to the opponent’s strategy, which implies that the outcome is a Nash equilibrium of the underlying game.

The following example demonstrates that even Nash equilibria may require distorted beliefs to be supported as BBE outcomes. Specifically, Example 1 shows that this is the case for Nash equilibrium in a Cournot competition. The intuition behind Example 1 is straightforward. The Cournot equilibrium cannot be supported by undistorted beliefs because such pairs of beliefs will induce one of the players to adopt a distorted belief by which he expects his opponent not to produce at all, and to best-reply to this distorted belief by producing the monopoly quantity. This in turn will force the opponent to reduce his production substantially below the Cournot level, making the deviator better off.

**Example 1** (*Cournot equilibrium cannot be supported by undistorted beliefs, yet it can be supported by blind beliefs*). Consider the following symmetric Cournot game  $G = (S, \pi)$ :  $S_i = [0, 1]$  and  $\pi_i(s_i, s_j) = s_i \cdot (1 - s_i - s_j)$  for each player  $i$ . The interpretation of the game is as follows. Each  $s_i$  is interpreted as the quantity chosen by firm  $i$ , the price of both goods is determined by the linear inverse demand function  $p = 1 - s_i - s_j$ , and the marginal cost of each firm is normalized to be zero. The unique Nash equilibrium of the game is  $s_i^* = s_j^* = \frac{1}{3}$ , which yields a payoff of  $\frac{1}{9}$  to both players. Assume to the contrary that this outcome can be supported as a weak BBE by the undistorted beliefs  $\psi_i^* = \psi_j^* = I_d$ . Consider a deviation of player 1 to the blind belief  $\psi_1' \equiv 0$ . The unique equilibrium of the biased game  $G_{(0, I_d)}$  is  $s_1' = \frac{1}{2}$ ,  $s_2' = \frac{1}{4}$ , which yields a payoff of  $\frac{1}{8} > \frac{1}{9}$  to the deviator. The unique Nash equilibrium  $s_i^* = s_j^* = \frac{1}{3}$  can be supported as the outcome of the strong BBE  $\left(\left(\frac{1}{3}, \frac{1}{3}\right), \left(\frac{1}{3}, \frac{1}{3}\right)\right)$  with blind beliefs, in which each agent believes the opponent is playing  $\frac{1}{3}$  regardless of the opponent’s actual play, and the agent plays the unique best reply to this belief, which is the strategy  $\frac{1}{3}$ .



*Remark 2* (Interpretation of Nash equilibria supported by blind beliefs.). We interpret an undistorted belief as describing an agent who has an accurate belief about the opponent's behavior on the equilibrium path, and, in addition, the agent keeps looking for cues that his opponent might have a different type, and if the agent observes such a cue, the agent evaluates the opponent's likely behavior, and best-responds to this assessment. Example 1 shows that the Cournot equilibrium cannot be supported by a population in which each agent keeps looking for cues for his opponent's type. In such a population, deviators would strictly earn by having a blind biased belief that induces the deviator to play the Stackelberg strategy. The incumbents will identify the mutants' type, and they will respond by playing the Stackelberg follower action, which will benefit the deviators.

By contrast, the second part of Example 1 (and its generalization in Proposition 1 below) shows that any Nash equilibrium can be supported by a blind belief, which is accurate on the equilibrium path. We interpret such a belief as describing an agent who understands correctly the equilibrium behavior of the opposing player, and ignores signals that suggest that his opponent is about to do something else. Our observation that it is rather equilibrium that supports belief rigidity, a prevalent behavioral phenomenon, and not disequilibrium is, we believe, quite interesting.

## 4.2 Any Nash Equilibrium is a BBE Outcome

The following result generalizes the second part of Example 1, and shows that any (strict) Nash equilibrium is an outcome of a (strong) BBE in which both players have blind beliefs that are accurate on the equilibrium path.

**Proposition 1.** *Let  $(s_1^*, s_2^*)$  be a (strict) Nash equilibrium of the game  $G = (S, \pi)$ . Let  $\psi_1^* \equiv s_2^*$  and  $\psi_2^* \equiv s_1^*$ . Then  $((\psi_1^*, \psi_2^*), (s_1^*, s_2^*))$  is a (strong) BBE.*

*Proof.* The fact that  $(s_1^*, s_2^*)$  is a Nash equilibrium of the underlying game implies that  $(s_1^*, s_2^*)$  is an equilibrium of the biased game  $G_{(\psi_1^*, \psi_2^*)}$ . The fact that the beliefs are blind implies that for any biased belief  $\psi'_i$ , there is an equilibrium in the biased game  $G_{(\psi'_i, \psi_j^*)}$  in which player  $j$  plays  $s_j^*$  and player  $i$  gains at most  $\pi_i(s_i^*, s_j^*)$ , which implies that  $((\psi_1^*, \psi_2^*), (s_1^*, s_2^*))$  is a BBE. Moreover, if  $(s_1^*, s_2^*)$  is a strict equilibrium, then in any equilibrium of any biased game  $G_{(\psi'_i, \psi_j^*)}$ , player  $j$  plays  $s_j^*$  and player  $i$  gains at most  $\pi_i(s_i^*, s_j^*)$ , which implies that  $((\psi_1^*, \psi_2^*), (s_1^*, s_2^*))$  is a strong BBE.  $\square$

An immediate corollary of Proposition 1 is that every game admits a BBE.

**Corollary 1.** *Every game admits a BBE.*

## 4.3 Zero-Sum Games

Recall that a game is *zero sum* if there exists  $c \in \mathbb{R}^+$  such that  $\pi_i(s_i, s_j) + \pi_j(s_i, s_j) = c$  for each strategy profile  $(s_i, s_j) \in S$ .

The following simple result shows that the unique Nash equilibrium payoff of a zero-sum game is also the unique payoff in any weak BBE.

*Claim 1.* The unique Nash equilibrium payoff of a zero-sum game is also the unique payoff in any weak BBE.

*Proof.* Let  $v_i$  be the unique Nash equilibrium payoff of player  $i$  in the underlying zero-sum game. Assume to the contrary that there exists a weak BBE  $(\psi^*, s^*)$  in which the payoff of player  $i$  is strictly lower than  $v_i$ . Consider a deviation of player  $i$  into the undistorted bias function  $\psi'_i = I_d$ . The assumption that  $(\psi^*, s^*)$  is a weak BBE implies that the deviator gets strictly less than  $v_i$  in a Nash equilibrium  $(s'_i, s'_j) \in NE\left(G_{(\psi'_i, \psi^*_j)}\right)$ , but this is impossible as the definition of  $v_i$  implies that there exists  $\hat{s}_i$  satisfying  $\pi_i(\hat{s}_i, s'_j) \geq v_i > \pi_i(s'_i, s'_j)$ .  $\square$

Example 6 in Appendix A.2 shows that even though the weak BBE payoff must be the Nash equilibrium payoff in a zero-sum game, the strategy profile sustaining it need not be a Nash equilibrium.

#### 4.4 Games with a Dominant Strategy

Next we show that if at least one of the players has a dominant strategy, then any weak BBE outcome must be a Nash equilibrium. Formally:

**Proposition 2.** *If a game admits a strictly dominant strategy  $s_i^*$  for player  $i$ , then any weak BBE outcome is a Nash equilibrium of the underlying game.*

*Proof.* Observe that  $s_i^*$  is the unique best reply of player  $i$  to any perceived strategy of player  $j$ , and, as a result, player  $i$  plays the dominant action  $s_i^*$  in any weak BBE. Assume to the contrary that there is a weak BBE in which player  $j$  does not best-reply against  $s_i^*$ . Consider a deviation of player  $j$  to choosing the undistorted belief  $I_d$ . Observe that player  $i$  still plays his dominant action  $s_i^*$ , and that player  $j$  best-responds to  $s_i^*$  in any Nash equilibrium of the induced biased game, and, as a result, player  $j$  achieves a strictly higher payoff, and we get a contradiction.  $\square$

Proposition 2 implies, in particular, that defection is the unique weak BBE outcome in the prisoner's dilemma game. Example 7 in Appendix A.1 demonstrates that a relatively small change to the prisoner's dilemma game, namely, adding a third weakly dominated “withdrawal” strategy that transforms “cooperation” into a weakly dominated strategy, allows us to sustain cooperation as a strong BBE outcome.

## 5 Monotone Games and Wishful Thinking

In this section we present a large class of games with monotone externalities and monotone differences, and define the notions of wishful thinking and pessimism, which will be analyzed in Section 6.

### 5.1 Monotone Games

We say that an interval game is monotone if it satisfies two conditions:

1. *Monotone externalities:* the payoff function of each player is strictly monotone in the opponent's strategy. Without loss of generality, we assume that the *externalities* are *positive*, i.e.,

the payoff of each player is increasing in the opponent's strategy, i.e., that  $\frac{\partial \pi_i(s_i, s_j)}{\partial s_j} > 0$  for each player  $i$  and each pair of strategies  $s_i, s_j$ . The assumption of positive externalities (given monotone externalities) is indeed without loss of generality because if originally the externalities with respect to player  $j$  are negative, then we can redefine player  $j$ 's strategy to be its inverse, and obtain positive externalities; for example, defining the difference between maximal capacity and quantity to be the strategy of each player in a Cournot competition yields a game with positive externalities.

In a game with positive externalities we refer to a player's strategy as his *investment*, and when  $s_i$  increases we refer to this increase as a larger investment by player  $i$ .

2. *Monotone differences*: For each player  $i$ , the derivative of the player's payoff with respect to his own strategy (i.e.,  $\frac{\partial \pi_i(s_i, s_j)}{\partial s_i}$ ) is strictly monotone in the opponent's strategy. Specifically, we divide the set of monotone games into three disjoint and exhaustive subsets:

- (a) *Strategic complements (increasing differences, supermodular games)*:  $\frac{\partial \pi_i(s_i, s_j)}{\partial s_i}$  is strictly increasing in  $s_j$  for each player  $i$  and each strategy  $s_i$  (or, equivalently,  $\frac{\partial^2 \pi_i(s_i, s_j)}{\partial s_i \partial s_j} > 0$  for each  $s_i, s_j$ ). Games with strategic complements are common in the economics literature, and include, in particular, price competitions with differentiated goods (Example 2 in Section ), input games (Example 9 in Appendix A.4), and stag-hunt games (Example 10 in Appendix A.4). Finite games with a payoff structure that resembles a discrete variant of strategic complements include the traveler's dilemma (Example 11 in Appendix A.4).
- (b) *Strategic substitutes (decreasing differences, submodular games)*:  $\frac{\partial \pi_i(s_i, s_j)}{\partial s_i}$  is strictly decreasing in  $s_j$  for each player  $i$  and each strategy  $s_i$  (or, equivalently,  $\frac{\partial^2 \pi_i(s_i, s_j)}{\partial s_i \partial s_j} < 0$  for each  $s_i, s_j$ ). Games with strategic substitutes are common in the economics literature, and include, in particular, Cournot (quantity) competitions (Example 3 below) and hawk-dove games (see Example 12 in Appendix A.5).
- (c) *Opposing differences*:  $\frac{\partial \pi_i(s_i, s_j)}{\partial s_i}$  is decreasing in  $s_j$  (for each strategy  $s_i$ ), while  $\frac{\partial \pi_j(s_i, s_j)}{\partial s_j}$  is increasing in  $s_i$  (for each strategy  $s_j$ ). Games with opposing differences are less common in the economics literature. Examples of these games include (1) duopolies in which one firm chooses its quantity, while the other firm chooses its price (see, e.g., Singh and Vives, 1984), and (2) asymmetric contests, in which it is often the case that a commitment of the favorite (underdog) player to exert more (less) effort induces the opponent to exert less effort (see, e.g., Dixit, 1987).

## 5.2 Wishful Thinking

We say that player  $i$  exhibits wishful thinking if the perceived opponent's strategy yields a higher payoff to the player relative to the real strategy the opponent plays. Formally:

**Definition 6.** Player  $i$  exhibits wishful thinking in weak BBE  $((\psi_1^*, \psi_2^*), (s_1^*, s_2^*))$  if  $\pi_i(s_i, \psi_i^*(s_j)) \geq \pi_i(s_i, s_j^*)$  for each  $s_i \in S_i$ .

*Remark 3.* Note that in a game with positive externalities player  $i$  exhibits wishful thinking in weak BBE  $((\psi_1^*, \psi_2^*), (s_1^*, s_2^*))$  iff  $\psi_2^*(s_1^*) \geq s_1^*$  and  $\psi_1^*(s_2^*) \geq s_2^*$ .

Similarly, we define the opposite notion, that of exhibiting pessimism. We say that a BBE exhibits pessimism if the perceived opponent's strategy yields a lower payoff to the player relative to the real opponent's strategy for all strategy profiles. It exhibits pessimism in equilibrium if it satisfies this property with respect to the strategy the opponent plays on the equilibrium path. Formally:

**Definition 7.** A weak BBE  $((\psi_1^*, \psi_2^*), (s_1^*, s_2^*))$  exhibits pessimism if  $\pi_i(s_i, \psi_i^*(s_j)) \leq \pi_i(s_i, s_j^*)$  for all  $s_i \in S_i$ .

### 5.3 Additional Definitions

In what follows we present two definitions that will be used in the analysis in the following sections: undominated Pareto optimality, and biased-belief minmax payoff.

We say that a strategy profile is undominated Pareto optimal if it is (1) undominated, and (2) Pareto optimal among all undominated strategy profiles. Formally:

**Definition 8.** Strategy profile  $(s_1^*, s_2^*)$  is *undominated Pareto optimal* if (1)  $s_i^* \in S_i^U$  for each player  $i$ , and (2) there does not exist  $(s'_1, s'_2) \in S_1^U \times S_2^U$  with a payoff that Pareto dominates  $(s_1^*, s_2^*)$ , i.e.,  $\pi_1(s_1^*, s_2^*) \leq \pi_1(s'_1, s'_2)$  and  $\pi_2(s_1^*, s_2^*) \leq \pi_2(s'_1, s'_2)$  where at least one of these inequalities is strict.

A biased-belief minmax payoff for player  $i$  (denoted by  $\tilde{M}_i^U$ ) is the maximal payoff player  $i$  can guarantee to himself in the following process: (1) player  $j$  chooses an arbitrary perceived strategy of player  $i$ , and (2) player  $i$  chooses a strategy profile, under the constraint that player  $j$ 's strategy is a best reply to the perceived strategy chosen above. That is,  $\tilde{M}_i^U$  is the payoff player  $i$  can guarantee himself no matter how his opponent (player  $j$ , she) perceives player  $i$ 's action, assuming that player  $j$  best-responds to what he believes player  $i$  is doing (and if there are multiple best replies, then we assume that player  $j$  chooses the best reply that is optimal for player  $i$ ). Formally:

**Definition 9.** Given game  $G = (A, u)$ , let  $\tilde{M}_i^U$ , the *biased-belief minmax payoff* of player  $i$ , be defined as follows:

$$\tilde{M}_i^U = \min_{s'_i \in S_i^U} \left( \max_{(s_i, s_j) \in S_i \times BR(s'_i)} \pi_i(s_i, s_j) \right).$$

Observe that the biased-belief minmax is weakly larger than the undominated maxmin (Definition 10), i.e.,  $\tilde{M}_i^U \geq M_i^U$  with an equality if the strategy of player  $j$  that guarantees that player  $i$ 's payoff is at most  $M_i^U$  is a unique best reply against some strategy of player  $i$  (which is the case, in particular, if the payoff function is strictly concave).

## 6 Main Results

Our main results characterize the set of BBE and BBE outcomes in three classes of games: games with strategic complements, games with strategic substitutes, and games with strategic opposites.

## 6.1 Preliminary Result: Necessary Conditions for a Weak BBE Outcome

We begin by defining undominated strategies and the undominated minmax payoff, which will be used to characterize necessary conditions for a strategy profile to be a weak BBE outcome.

Strategy  $s_i$  of player  $i$  is *undominated* if it is a best reply of some strategy of the opponent, i.e., if there exists strategy  $s_j \in S_j$ , such that  $s_i \in BR(s_j)$ . We say that a strategy profile is *undominated* if both strategies in the profile are undominated. Recall that in a finite game, due to the minmax theorem, a strategy is undominated iff it is not strictly dominated by another strategy.

Let  $S_i^U \in S_i$  denote the *set of undominated strategies* of player  $i$ . Observe that  $S_i^U$  is not necessarily a convex set.

An undominated minmax payoff for player  $i$  is the maximal payoff player  $i$  can guarantee to himself in the following process: (1) player  $j$  chooses an arbitrary undominated strategy, and (2) player  $i$  chooses a strategy (after observing player  $j$ 's strategy). Formally:

**Definition 10.** Given game  $G = (S, u)$ , let  $M_i^U$ , the *undominated minmax payoff* of player  $i$ , be defined as follows:

$$M_i^U = \min_{s_j \in S_j^U} \left( \max_{s_i \in S_i} \pi_i(s_i, s_j) \right).$$

Observe that the undominated minmax is weakly larger than the standard maxmin, i.e.,  $M_i^U \geq \min_{s_j \in S_j} (\max_{s_i \in S_i} \pi_i(s_i, s_j))$  with an equality if player  $j$  does not have any strictly dominated strategy<sup>2</sup> (i.e., if  $S_j^U = S_j$ ).

The following simple result (which will be helpful in deriving the main results in the following subsections) shows that any weak BBE outcome is an undominated strategy profile that yields a payoff above the player's undominated minmax payoff to each player.

**Proposition 3.** *If a strategy profile  $s^* = (s_1^*, s_2^*)$  is a weak BBE outcome, then (1) the profile  $s^*$  is undominated and (2)  $\pi_i(s^*) \geq M_i^U$ .*

*Proof.* Assume that  $s^* = (s_1^*, s_2^*)$  is a biased-belief equilibrium outcome. This implies that each  $s_i^*$  is a best reply to the player's distorted belief, which implies that each  $s_i^*$  is undominated. Assume to the contrary, that  $\pi_i(s^*) < M_i^U$ . Then, by deviating to the undistorted function  $I_d$ , player  $i$  can guarantee a fitness of at least  $M_i^U$  in any distorted equilibrium.  $\square$

## 6.2 Games with Strategic Complements

Our first main result characterizes the set of BBE outcomes in games with strategic complements. It shows that a strategy profile is a BBE outcome essentially iff (I) it is undominated, (II) it yields a payoff above the undominated/biased-belief minmax payoff to both players, and (III) both players overinvest (i.e., use a weakly higher strategy than the best reply to the opponent). Formally:

**Proposition 4.** *Let  $G$  be a game with strategic complements and positive externalities.*

---

<sup>2</sup>The undominated minmax payoff might be strictly higher than the undominated *maxmin* payoff due to the non-convexity of  $S_j^U$ ; i.e., player  $i$  might be able to guarantee only a lower payoff in a setup in which player  $j$  is allowed to choose his undominated strategy after observing player  $i$ 's chosen strategy.

1. Let  $(s_1^*, s_2^*)$  be a BBE outcome. Then  $(s_1^*, s_2^*)$  has the following properties: (I) it is undominated, and it satisfies for each player  $i$ : (II)  $\pi_i(s_i^*, s_j^*) \geq M_i^U$ , and (III) overinvestment:  $s_i^* \geq \min(BR(s_j^*))$ .
2. Let  $(s_1^*, s_2^*)$  be an undominated profile that satisfies, for each player  $i$ : (II)  $\pi_i(s_i^*, s_j^*) > \tilde{M}_i^U$ , and (III)  $s_i^* \geq \min(BR(s_j^*))$ . Then,  $(s_1^*, s_2^*)$  is a BBE outcome.  
Moreover, if  $\pi_i(s_i, s_j)$  is strictly concave in  $s_i$  (i.e.,  $\frac{\partial^2 \pi_i(s_i, s_j)}{\partial s_i^2} > 0$ ) then  $(s_1^*, s_2^*)$  is a strong BBE outcome.

*Sketch of Proof (formal proof in Appendix F.2).*

**Part 1:** Proposition 3 implies (I) and (II). To prove (III, overinvestment), assume to the contrary that  $s_i^* < \min(BR(s_j^*))$ . Consider a deviation of player  $i$  that induces him to invest slightly more than  $s_i^*$ . The fact that  $s_i^* < \min(BR(s_j^*))$  implies that player  $i$  strictly earns from his own deviation. The assumption that the biased belief of the opponent is monotone implies that the agent's deviation induces the opponent to invest more and, thereby to further improve the agent's payoff. Thus, the agent gains from the deviation, and  $(s_1^*, s_2^*)$  cannot be a BBE outcome.

**Part 2:** The strategy profile  $(s_1^*, s_2^*)$  is supported as a BBE outcome by a profile of biased beliefs  $(\psi_1^*, \psi_2^*)$  in which each biased belief  $\psi_j^*$  satisfies: (1) blindness to good news:  $\psi_j^*$  distorts any  $s'_i \geq s_i^*$  into  $BR^{-1}(s_j^*)$ , and (2) overreaction to bad news:  $\psi_j^*$  distorts any  $s'_i < s_i^*$  to a sufficiently low strategy  $\psi_j(s'_i)$ , such that player  $i$  loses in any strategy profile  $(s'_i, s'_j)$  in which player  $j$  best-responds to the perceived strategy of player  $i$  (i.e.,  $s'_j \in BR(\psi_j(s'_i))$ ). These properties imply that  $((\psi_1^*, \psi_2^*), (s_1^*, s_2^*))$  is a BBE (and a strong BBE if the payoff function is strictly concave).  $\square$

Recall that a game with strategic complements admits a *lowest Nash equilibrium*  $(\underline{s}_1, \underline{s}_2)$  in which both players invest less than in any other Nash equilibrium, i.e.,  $s'_i \geq \underline{s}_i$  for each player  $i$  and each strategy  $s'_i$  that is played in a Nash equilibrium (see, e.g., Milgrom and Roberts, 1990).

An immediate corollary of Prop. 4 is that in each BBE outcome, both players invest more than in any Nash equilibrium. Formally:

**Corollary 2.** *Let  $G$  be a game with strategic complements and positive externalities with a lowest Nash equilibrium  $(\underline{s}_1, \underline{s}_2)$  that satisfies  $\underline{s}_1 < \max(S_i)$  for each player  $i$ . Let  $(s_1^*, s_2^*)$  be a BBE outcome. Then  $\underline{s}_i \leq s_i^*$  for each player  $i$ .*

*Proof.* The result is immediate from part (1.III) of Proposition 4 (namely, that both agents weakly overinvest in any BBE outcome), and the observation (which is formally proved in Lemma 1 in Appendix F.3) that  $s_i^* < \underline{s}_i$  implies that at least one of the players strictly underinvests.  $\square$

Corollary 2 shows that the notion of BBE rules out socially bad outcomes in which one (or both) of the players invests less effort than the lowest Nash equilibrium. In particular, in a price competition with differentiated goods (see Example 2 below), the corollary implies that the price chosen by any player in any BBE is at least the player's price in the unique Nash equilibrium of the game.



The final corollary shows the close relation between BBE and wishful thinking. Specifically, it shows that any biased belief in any BBE (with a non-extreme outcome) of a game with strategic complements exhibits wishful thinking. The intuition is that wishful thinking causes an agent to believe that the opponent is playing a higher action, which induces the agent to respond with a higher action, which, in turn, causes the opponent to respond by playing a higher action, which benefits the agent.<sup>3</sup>

**Corollary 3.** *Let  $G$  be a game with positive externalities and strategic complements. Let  $((\psi_1^*, \psi_2^*), (s_1^*, s_2^*))$  be a BBE. If  $s_i^* \notin \{\min(S_i), \max(S_i)\}$ , then player  $i$  exhibits wishful thinking (i.e.,  $\psi_i^*(s_j^*) \geq s_j^*$ ).*

*Proof.* Assume to the contrary that  $\psi_i^*(s_j^*) < s_j^*$ . The strategic complementarity implies that  $\max(BR(\psi_i^*(s_j^*))) \leq \min(BR(s_j^*))$  with an equality only if

$$\max(BR(\psi_i^*(s_j^*))) \in \{\min(S_i), \max(S_i)\}$$

(see Lemma 2 in Appendix F.4 for a formal proof of this claim). Part 1 of Proposition 4 and the definition of a BBE imply that

$$\max(BR(\psi_i^*(s_j^*))) \geq s_i^* \geq \min(BR(s_j^*)).$$

The previous inequalities jointly imply that

$$\max(BR(\psi_i^*(s_j^*))) = s_i^* = \min(BR(s_j^*)) \in \{\min(S_i), \max(S_i)\},$$

which contradicts the assumption that  $s_i^* \notin \{\min(S_i), \max(S_i)\}$ .  $\square$

Next, we apply our analysis of games with strategic complements to price competition with differentiated goods (the linear city model à la Hotelling). Specifically, we show that (1) players choose prices above the unique Nash equilibrium price in all BBE, and (2) any undominated symmetric price profile above the Nash equilibrium price can be supported as a strong BBE. In Appendix A.4 we present three additional examples: input games, stag hunt games, and the traveler's dilemma.

**Example 2** (*Price competition with differentiated goods; see a textbook analysis in Mas-Colell, Whinston, and Green, 1995, Section 12.C*). Consider a mass one of consumers equally distributed in the interval  $[0, 1]$ . Consider two firms that produce widgets, located at the two extreme locations: 0 and 1. Every consumer wants at most one widget. Producing a widget has a constant marginal cost, which we normalize to be zero. Each firm  $i$  chooses price  $s_i \in [0, M]$  for its widgets. The total cost of buying a widget from firm  $i$  is equal to its price  $s_i$  plus  $t$  times the consumer's distance from

<sup>3</sup>Corollary 3 allows for pessimism of player  $i$  in a BBE only if player  $i$  plays an extreme strategy (either, the minimal feasible strategy or the maximal feasible strategy) and his pessimism does not affect his play; i.e., the best reply against the real opponent's strategy and the best reply against the perceived opponent's strategy coincide in being the same extreme strategy. For example, this is the case in the biased beliefs that support the action profile  $(s_i, s_j)$  in the stag hunt game analyzed below.

the firm, where  $t \in [0, M]$ ). Each buyer buys a widget from the firm with the lower total buying cost. This implies that the total demand for good  $i$  is given by function  $q_i(s_i, s_j)$ :

$$q_i(s_i, s_j) = \begin{cases} 0 & \frac{s_j - s_i + t}{2 \cdot t} < 0 \\ \frac{s_j - s_i + t}{2 \cdot t} & 0 < \frac{s_j - s_i + t}{2 \cdot t} < 1 \\ 1 & \frac{s_j - s_i + t}{2 \cdot t} > 1, \end{cases}$$

The payoff (profit) of firm  $i$  is given by  $\pi_i(s_i, s_j) = s_i \cdot q_i(s_i, s_j)$ . Observe that the payoff function is strictly concave in  $s_i$  for any non-extreme  $s_j$  (and it is weakly concave for the extreme values of  $s_j$ ). One can show that the game has strategic complements, and that the best-reply function of each player is:

$$s_i(s_j) = \begin{cases} \frac{s_j + t}{2} & s_j < 3 \cdot t \\ s_j - t & s_j \geq 3 \cdot t. \end{cases}$$

It is well known that the unique Nash equilibrium of this example is given by  $s_i = s_j = t$ , which yields a payoff of  $\frac{t}{2}$  to each firm.

Observe that the set of undominated strategies of each player is the interval  $\left[\frac{t}{2}, \frac{M+t}{2}\right]$  (where  $\frac{t}{2}$  is the best reply against 0 and  $\frac{M+t}{2}$  is the best reply against  $M$ ). This implies that the undominated minmax of each player is equal to  $\pi_i\left(\frac{3}{4} \cdot t, \frac{t}{2}\right) = \frac{3}{4} \cdot t \cdot \frac{3}{8} = \frac{9}{32} \cdot t$ . Proposition 4 implies that a strategy profile  $(s_i, s_j)$  is a BBE outcome if for each player  $i$ : (1)  $s_i \in \left[\frac{t}{2}, \frac{M+t}{2}\right]$  (undominated strategy), (2)  $\pi_i(s_i, s_j) > \frac{9}{32} \cdot t$  (payoff above the undominated minmax payoff),<sup>4</sup> and (3) overinvestment:  $s_i \geq \frac{s_j + t}{2}$ .

Figure 1 shows the set of BBE outcomes (which coincides with the set of strong BBE outcomes, due to the strict concavity of the payoff function), for  $t = 1$  and  $M = 3$ .

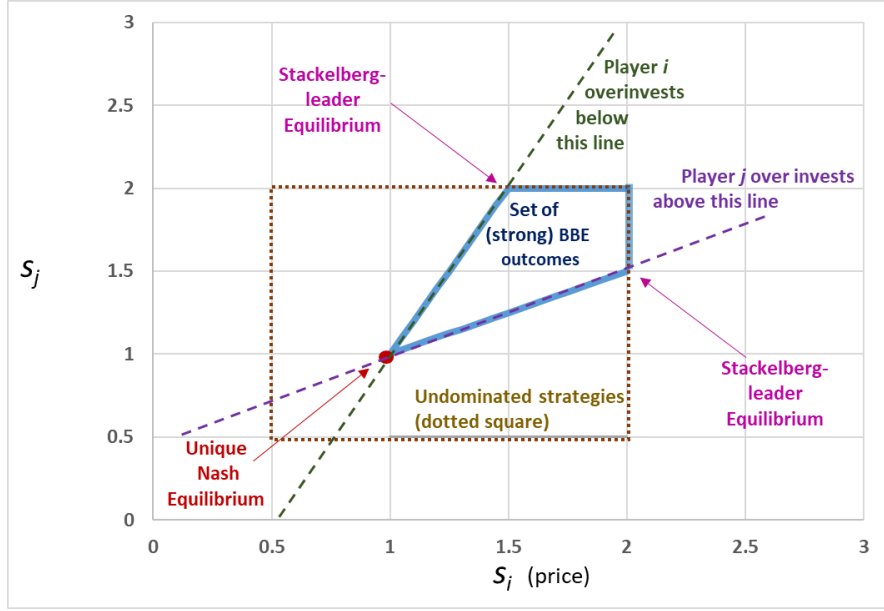
Observe that the sum of the payoffs to the two firms,  $s_i \cdot q_i(s_i, s_j) + s_j \cdot q_j(s_i, s_j)$ , is a mixed average of  $s_i$  and  $s_j$ . The fact that the Nash equilibrium is in the bottom left corner of the set of BBE outcomes implies that all BBE outcomes (except the Nash equilibrium itself) strictly improve social welfare relative to the Nash equilibrium (as measured by the sum of payoffs of the two firms).

Next, we make two observations regarding the implications of the extent of wishful thinking on the players' payoffs (both observations hold also for the input games in Example 9 in Appendix A.4):

1. Increasing the wishful thinking of both players improves the players' payoffs. Specifically, with respect to symmetric BBE outcomes, a higher level of wishful thinking induces a higher equilibrium price and a higher payoff to the players: a wishful thinking level of  $x^* \equiv \psi^*(s^*) - s^* \in [0, 1]$  induces the symmetric BBE price  $x^* + t = x^* + 1$  (which is implied by the perceived best-reply equation  $s^* = \frac{\psi^*(s^*) + t}{2} = \frac{s^* + x^* + t}{2}$ ), which yields a payoff of  $\frac{x^* + 1}{2}$  to each player.

<sup>4</sup>One can show that the constraint on  $s_i$  implied by  $\pi_i(s_i, s_j) > \frac{9}{32} \cdot t$  is nonbinding. The constraint is

$$s_i \in \left( \frac{s_j + t - \sqrt{(s_j + t)^2 - 2.25 \cdot t^2}}{2}, \frac{s_j + t - \sqrt{(s_j + t)^2 - 2.25 \cdot t^2}}{2} \right).$$

Figure 1: The Set of (Strong) BBE Outcomes in Example 2 ( $t = 1, M = 3$ )

2. When the wishful thinking levels of the two players differ, the player with the higher wishful thinking level has a lower payoff. This is because the difference between the payoffs of a firm with price  $s_i$  and an opponent with price  $s_j < s_i$  is equal to:

$$\pi_i - \pi_j = s_i \cdot \left( \frac{s_j - s_i + 1}{2} \right) - s_j \cdot \left( \frac{s_i - s_j + 1}{2} \right) = 0.5 (s_j (s_j - 1) - s_i (s_i - 1)) < 0.$$

Intuitively, wishful thinking is like a public good in this setup: (1) a higher level of wishful thinking is beneficial to social welfare, and (2) if the two players have different levels of wishful thinking, the player with the higher level obtains a lower payoff.

We conclude the example by presenting a symmetric biased belief  $\psi_1^* = \psi_2^*$  that supports the outcome  $(2, 2)$  as the BBE  $((\psi_1^*, \psi_2^*), (2, 2))$  in the game with  $M = 3$  and  $t = 1$ :

$$\psi_i^*(s_j) = \begin{cases} 3 & s_j > 2 \\ 2 \cdot s_j - 1 & s_j \in [0.5, 2] \\ 0 & s_j < 0.5. \end{cases}$$

Observe that: (1) this BBE yields a payoff of 1 to each player and (2) the biased belief presents wishful thinking, i.e.,  $\psi_i^*(2) = 3 > 2$ . Further observe that a player with biased belief  $\psi_i^*$  plays the same strategy as the opponent (regardless of the opponent's biased belief) in any equilibrium of the biased game in which the opponent plays any intermediate value of  $s_j$  (i.e.,  $s_j \in [0.5, 2]$ ):

$$s_i(\psi_i^*(s_j)) = \begin{cases} s_i(3) = 2 & s_j > 2 \\ s_i(2 \cdot s_j - 1) = 0.5 \cdot (2 \cdot s_j - 1 + 1) = s_j & s_j \in [0.5, 2] \\ s_i(0) = 0.5 & s_j < 0.5. \end{cases}$$

This implies that the equilibrium payoff of a deviating player  $j$  who plays strategy  $s_j$  is equal to:

$$\pi_j(s_j, s_i(\psi_i^*(s_j))) = \begin{cases} s_j \cdot 0.5 \cdot (2 - s_j + 1) = s_j \cdot 0.5 \cdot (3 - s_j) < 1 & s_j > 2 \\ s_j \cdot 0.5 \cdot q(s_j, s_j) = 0.5 \cdot s_j & s_j \in [0.5, 2] \\ s_j \cdot 0.5 \cdot (0.5 - s_j + 1) = s_j \cdot 0.5 \cdot (1.5 - s_j) < 0.25 & s_j < 0.5, \end{cases}$$

and it is at most 1, which implies that a deviator cannot gain from his deviation.

Finally, note that Figure 1 shows that the two Stackelberg-leader equilibria (the unique subgame-perfect equilibrium of the sequential games in which one of the players plays first, and the opponent replies after observing the leader's strategy) are included in the set of BBE, as is proven in general in Proposition 7.

### 6.3 Games with Strategic Substitutes

Our next result characterizes the set of BBE outcomes in games with strategic substitutes (and positive externalities). It shows that a strategy profile is a BBE outcome essentially iff (I) it is undominated, (II) it yields a payoff above the undominated/biased-belief minmax payoff to both players, and (III) both players underinvest (i.e., use a weakly lower strategy than the best reply to the opponent). Formally:

**Proposition 5.** *Let  $G$  be a game with strategic substitutes and positive externalities.*

1. *Let  $(s_1^*, s_2^*)$  be a BBE outcome. Then  $(s_1^*, s_2^*)$  has the following properties: (I) it is undominated, and it satisfies for each player  $i$ : (II)  $\pi_i(s_i^*, s_j^*) \geq M_i^U$  and (III)  $s_i^* \leq \max(BR(s_j^*))$  (underinvestment).*
2. *Let  $(s_1^*, s_2^*)$  be an undominated profile that satisfies, for each player  $i$ : (II)  $\pi_i(s_i^*, s_j^*) > \tilde{M}_i^U$ , and (III)  $s_i^* \leq \max(BR(s_j^*))$ . Then,  $(s_1^*, s_2^*)$  is a BBE outcome. Moreover, if  $\pi_i(s_i, s_j)$  is strictly concave then  $(s_1^*, s_2^*)$  is a strong BBE outcome.*

The proof, which is analogous to the proof of Proposition 4, is presented in Appendix F.5.

An immediate corollary of Proposition 4 is that in each BBE outcome, at least one of the players invests less relative to his maximal Nash equilibrium investment. Formally:

**Corollary 4.** *Let  $G$  be a game with strategic substitutes and positive externalities. Let  $(s_1^*, s_2^*)$  be a BBE outcome. Then, there exists a Nash equilibrium of the underlying game  $(s_1^e, s_2^e)$ , and a player  $i$  such that  $s_i^e \geq s_i^*$ .*

*Proof.* The result is immediate from part (1-III) of Proposition 4 (namely, that both agents weakly underinvest in any BBE outcome), and the observation (which is formally proved in Lemma 1 in Appendix F.3) that if the effort of each player  $s_i^*$  is strictly below all of his Nash equilibrium efforts, then at least one of the players strictly underinvests.  $\square$

Corollary 4 shows that the notion of BBE rules out socially good outcomes in which both players invest more effort relative to their maximal Nash equilibrium effort. In particular, in a Cournot

competition (see Example 3 below), the corollary implies that a collusive outcome in which both players retain more unused capacity relative to the unique Nash equilibrium.

Combining Corollary 2 and Corollary 4 implies the following *empirical prediction of our model and the notion of BBE: efficient (non-Nash equilibrium) outcomes are easier to support in games with strategic complements, relative to games with strategic substitutes*. This prediction is consistent with the experimental findings of Potters and Suetens (2009), which show that there is significantly more cooperation in games with strategic complements than in games with strategic substitutes.

The following corollary shows that in games with strategic substitutes, as in games with strategic complements, there is a close relation between BBE and wishful thinking. Specifically, it shows that any biased belief in any BBE (with a non-extreme outcome) of a game with strategic substitutes exhibits wishful thinking. The intuition is that wishful thinking causes an agent to believe that the opponent is playing a higher action, which induces the agent to respond with a lower action, which, in turn, causes the opponent to respond by playing a higher action, which benefits the agent.

**Corollary 5.** *Let  $G$  be a game with positive externalities and strategic substitutes. Let  $((\psi_1^*, \psi_2^*), (s_1^*, s_2^*))$  be a BBE. If  $s_i^* \notin \{\min(S_i), \max(S_i)\}$ , then player  $i$  exhibits wishful thinking (i.e.,  $\psi_i^*(s_j^*) \geq s_j^*$ ).*

The proof, which is analogous to the proof of Corollary 3, is presented in Appendix F.7.

The following example characterizes the set of BBE outcomes in a Cournot competition. Appendix A.5 presents an analysis of another game of strategic substitutes: the hawk-dove game.

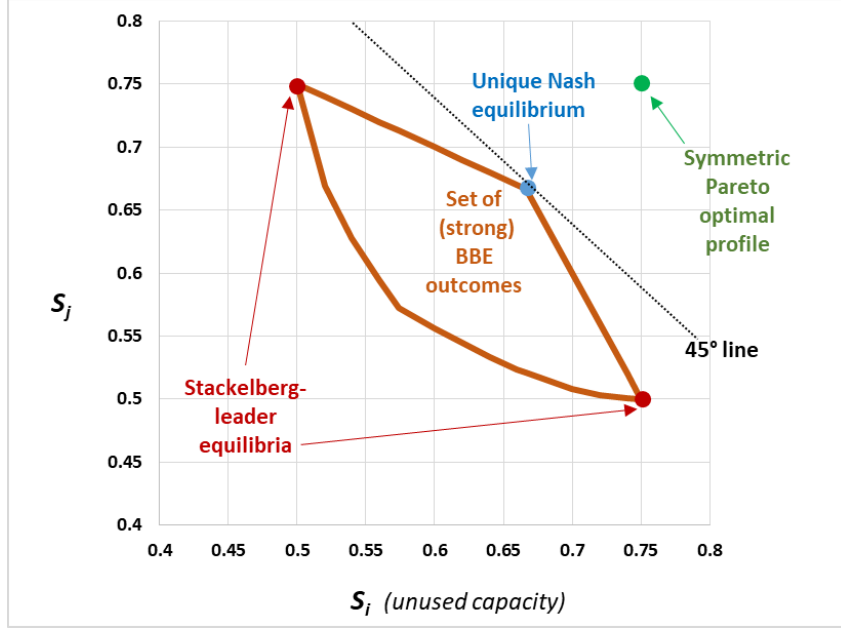
**Example 3** (*Cournot competition with linear demand*). Consider a symmetric Cournot competition, where we relabel the set of strategies to describe unused capacity, rather than quantity, in order to follow the normalization of positive externalities. Formally, let  $G = (S, \pi)$ :  $S_i = [0, 1]$  and  $\pi_i(s_i, s_j) = (1 - s_i) \cdot (s_i + s_j - 1)$  for each player  $i$ . Each  $s_i$  is interpreted as the unused capacity (= one minus the quantity, i.e.,  $s_i = 1 - q_i$ ) chosen by firm  $i$ , the price of both goods is determined by the linear inverse demand function  $p = 1 - q_i - q_j = s_i + s_j - 1$ , and the marginal cost of each firm is normalized to be zero.

Observe that:

1.  $BR(s_i) = 1 - \frac{s_i}{2}$ , and the unique Nash equilibrium of the game is  $s_1^* = s_2^* = \frac{2}{3}$ , which yields a payoff of  $\frac{1}{9}$  to both players.
2. The set of undominated strategies of each player is the interval  $[0.5, 1]$  (where 1 is the best reply against 0, and 0.5 is the best reply against 1).
3. The symmetric Pareto optimal profile (which is also undominated) is  $s_i = s_j = \frac{3}{4}$ , yielding a payoff of  $\frac{1}{8}$  to each player.
4. The undominated minmax payoff  $M_i^U = \frac{1}{16}$ , which is achieved by the opponent playing his lowest undominated strategy  $s_i = 0.5$ .
5. The sum of payoffs of both players when they play profile  $(s_1, s_2)$  is  $\pi_1(s_1, s_2) + \pi_2(s_1, s_2) = (2 - (s_i + s_j)) \cdot ((s_i + s_j) - 1)$ , which is an increasing function of  $s_i + s_j$  in the domain of undominated strategies  $s_i, s_j \geq 0.5$ .

Applying the analysis of the previous subsection to a Cournot competition shows that strategy profile  $(s_1, s_2)$  is a BBE outcome iff it satisfies for each player  $i$ : (1) the strategy is undominated:  $s_i \geq 0.5$ , (2) the payoff is greater than the undominated minmax payoff:  $(1 - s_i) \cdot (s_i + s_j - 1) \geq \frac{1}{16} = M_i^U$ , and (3) underinvestment relative to the best reply against the opponent:  $s_i \leq BR(s_j) = 1 - \frac{s_j}{2}$ . Due to having a strictly concave payoff function, the set of BBE outcomes coincides with the set of strong BBE outcomes. Figure 1 shows this set of BBE outcomes (the strategy profiles that satisfy the above three conditions).

Figure 2: The Set of (Strong) BBE Outcomes in a Cournot Competition



Observe that the unique Nash equilibrium  $(\frac{2}{3}, \frac{2}{3})$  is the profile that maximizes the sum  $s_i + s_j$  within the set of BBE. This implies that all other BBE outcomes yield lower social welfare (as measured by the sum of payoffs) relative to the Nash equilibrium.

Next, we make two observations regarding the implications of the level of wishful thinking on the players' payoffs.

1. Increasing the wishful thinking of both players decreases the players' payoffs. Specifically, when focusing on symmetric BBE outcomes, a higher level of wishful thinking induces a lower level of unused capacity and a lower payoff to both players; the higher level of production is induced by the false assessment of each firm that the other firm is producing less than it actually does.<sup>5</sup>
2. When the wishful thinking levels of the two players differ, the player with the higher wishful thinking has a higher payoff. This is because the difference between the payoffs of a firm with

<sup>5</sup>A wishful thinking level of  $x^* \equiv \psi^*(s^*) - s^* \in [0, 0.28]$  induces a symmetric BBE unused capacity of  $s^* = \frac{2-x^*}{3}$  (which is implied by the perceived best-reply equation  $s^* = 1 - \frac{\psi^*(s^*)}{2} = 1 - \frac{s^* + x^*}{2}$ ).



price  $s_i$  and an opponent with price  $s_j < s_i$  is equal to

$$\pi_i - \pi_j = s_i \cdot \left( \frac{s_j - s_i + 1}{2} \right) - s_j \cdot \left( \frac{s_i - s_j + 1}{2} \right) = 0.5 (s_j (s_j - 1) - s_i (s_i - 1)) < 0.$$

Thus, a higher level of wishful thinking is beneficial to social welfare, but harms the player with the higher level (relative to the opponent's payoff).

Finally, note that Figure 2 shows that the two Stackelberg-leader equilibria (the unique subgame-perfect equilibria of the sequential games in which one of the players plays first, and the opponent replies after observing the leader's strategy) are included in the set of BBE, as is proven in general in Proposition 7.

## 6.4 Pessimism in Games with Opposing Differences

The results of the previous two subsections present a strong tendency of BBE to exhibit wishful thinking both in games with strategic complements and in games with strategic substitutes. This raises the question of which class of games induces pessimism. In this section we show that the answer to this question is games with strategic opposites. Recall that these are games in which the strategy of player 1 is a complement of player 2's strategy, while the strategy of player 2 is a substitute of player 1's strategy, e.g., duopolistic competitions in which one firm chooses its quantity while the opposing firm chooses its price (Singh and Vives, 1984) and various classes of asymmetric contests (Dixit, 1987).

Proposition 6 characterizes the set of BBE outcomes in games with strategic opposites (and positive externalities).

It shows that a strategy profile is a BBE outcome essentially iff (I) it is undominated, (II) it yields a payoff above the undominated/biased-belief minmax payoff to both players, and (III) player 1 (for whom player 2's strategy is a complement) underinvests, while player 2 (for whom player 1's strategy is a substitute) overinvests. Formally:

**Proposition 6.** *Let  $G$  be a game with positive externalities and strategic opposites:  $\frac{\partial^2 \pi_1(s_1, s_2)}{\partial s_1 \partial s_2} > 0$  and  $\frac{\partial^1 \pi_2(s_1, s_2)}{\partial s_1 \partial s_2} < 0$  for each pair of strategies  $s_1, s_2$ .*

1. *Let  $(s_1^*, s_2^*)$  be a BBE outcome. Then  $(s_1^*, s_2^*)$  is (I) undominated: (II)  $\pi_i(s_i^*, s_j^*) \geq M_i^U$  for each player  $i$ , and (III)  $s_1^* \leq \max(BR(s_2^*))$  and  $s_2^* \geq \min(BR(s_1^*))$  (i.e., player 1 underinvests and player 2 overinvests relative to the best reply to the opponent).*
2. *Let  $(s_1^*, s_2^*)$  be a profile satisfying the following conditions: (I)  $(s_1^*, s_2^*)$  is undominated, (II)  $\pi_i(s_i^*, s_j^*) > \tilde{M}_i^U$  for each player  $i$ , and (III)  $s_1^* \leq \max(BR(s_2^*))$  and  $s_2^* \geq \min(BR(s_1^*))$ . Then,  $(s_1^*, s_2^*)$  is a BBE outcome.*

The proof, which is analogous to the proof of Proposition 4, is presented in Appendix F.8.

The following corollary shows that in games with strategic opposites, there is a close relation between BBE and pessimism. Specifically, it shows that any biased belief in any BBE (with a non-extreme outcome) of a game with strategic opposites exhibits pessimism. The intuition is that

pessimism causes player 1 to believe that player 2 is playing a lower action, which induces player 1 to respond with a lower action, which, in turn, causes player 2 to respond by playing a higher action, which benefits player 1. Similarly, pessimism causes player 2 to believe that player 1 is playing a lower action, which induces player 2 to respond with a higher action, which, in turn, causes player 1 to respond by playing a higher action, which benefits player 2.

**Corollary 6.** *Let  $((\psi_1^*, \psi_2^*), (s_1^*, s_2^*))$  be a BBE of a game with positive externalities and strategic opposites (i.e.,  $\frac{\partial^2 \pi_1(s_1, s_2)}{\partial s_1 \partial s_2} > 0$  and  $\frac{\partial^2 \pi_2(s_1, s_2)}{\partial s_1 \partial s_2} < 0$  for each pair of strategies  $s_1, s_2$ ). If  $s_i^* \notin \{\min(S_i), \max(S_i)\}$ , then player  $i$  exhibits pessimism (i.e.,  $\psi_i^*(s_j^*) \leq s_j^*$ ).*

The proof, which is analogous to the proof of Corollary 3, is presented in Appendix F.9.

Next, we present an example of a game with strategic opposites, and we characterize the set of BBE in this game.

**Example 4** (*Matching pennies with positive externalities*). The game presented in Table 1, a variant of the matching pennies game, is played as follows:

1. Player 1 (player 2) gains 1 utility point from matching (mismatching) his opponent.
2. Each player  $i$  induces a gain of 3 utility points to his opponent by choosing heads (action  $h_i$ ).

The game admits a unique Nash equilibrium  $(0.5, 0.5)$  with a payoff of 1.5 to each player. The (undominated) minmax payoff of each player is 1 (obtained when the opponent plays  $t_j$ ). Observe that the game has positive externalities, that the strategy of player 2 is a strategic complement for player 1, while the strategy of player 1 is a strategic substitute for player 2.

Table 1: Matching Pennies with Positive Externalities

	$h_2$	$t_2$
$h_1$	4, 2	-1, 4
$t_1$	2, 1	1, -1

Applying the analysis of the previous section shows that the game admits 2 classes of BBE:

1. A class in which the players mix while giving a larger weight to playing heads (the action with positive externalities), pessimism, and one-directional blindness. Specifically, each BBE in this class  $((\psi_1^*, \psi_2^*), (\beta_1, \beta_2))$  satisfies for each player  $i$ : (I)  $\beta_i \in [0.5, 1]$  (i.e., both players play heads more frequently than in the unique Nash equilibrium), (II) pessimism:  $\psi_i^*(\beta_j) = 0.5 < \beta_j$ , and (III) one-sided blindness:  $\psi_1^*(\alpha) = 0.5$  for each  $\alpha \geq \beta_2$ ;  $\psi_1^*(\alpha) < 0.5$  for each  $\alpha < \beta_1$ ;  $\psi_2^*(\alpha) = 0.5$  for each  $\alpha \leq \beta_2$ ; and  $\psi_2^*(\alpha) < 0.5$  for each  $\alpha > \beta_2$ .
2. A class in which player 1 mixes while giving more weight to tails, while player 2 plays heads. Both players exhibit pessimism. Specifically, each BBE in this class  $((\psi_1^*, \psi_2^*), (\beta_1, \beta_2))$  satisfies for each player  $i$ : (I)  $\beta_1 \in [0, 0.5]$  and  $\beta_2 = 1$  (i.e., player 1 plays tails more frequently than in the unique Nash equilibrium, while player 2 always plays heads), (II) pessimism for player 1:

$\psi_1^*(\beta_2 = 1) = 0.5 < 1$ , and  $\psi_2^*(\beta_1) = 0.5$  (player 2 is not pessimistic, due to the fact that he chooses the extreme action 1), and (III)  $\psi_2^*(\alpha) > 0.5$  for each  $\alpha > \beta_1$ .

Observe that any profile  $(\beta_1, \beta_2)$ , where  $\beta_2 < 0.5$  or  $(\beta_1 < 0.5 \text{ and } \beta_2 < 1)$ , cannot be a BBE outcome:

1. If  $\beta_2 < 0.5$  and  $\beta_1 = 0$ , then player 2's payoff is negative, and less than his undominated minmax payoff of 1.
2. If  $\beta_2 < 0.5$  and  $\beta_1 > 0$ , then player 1 can gain by deviating to  $\psi'_1 \equiv 0$ , as the only possible equilibria of the new biased game are  $(0_1, 0_2)$  and  $(0_1, \beta_2)$ , both of which induce a higher payoff to player 1 relative to  $(\beta_1, \beta_2)$ .
3. If  $\beta_1 < 0.5$  and  $\beta_2 < 1$ , then player 2 can gain by deviating to  $\psi'_2 \equiv 0$ , as the only possible equilibria of the new biased game are  $(0_1, 1_2)$  and  $(\beta_1, 1_2)$ , both of which induce a higher payoff to player 2 relative to  $(\beta_1, \beta_2)$ .

## 6.5 Empirical Prediction Regarding Wishful Thinking

Arguably, the class of games with strategic opposites (which induces pessimism) is less common in strategic interactions than the classes of games with strategic complements/substitutes (both of which induce wishful thinking). This observation suggests the following empirical predictions of our model: (1) wishful thinking is more common than pessimism, and (2) there are some (less common) strategic interactions that induce pessimism. This empirical prediction is consistent with the experimental evidence that people tend to present wishful thinking, although, the extent of wishful thinking may substantially differ across different environments and may disappear in some environments (see, e.g., [Babad and Katz, 1991](#); [Budescu and Bruderman, 1995](#); [Bar-Hillel and Budescu, 1995](#); [Mayraz, 2013](#)).

## 7 Additional Results

### 7.1 BBE with Strategic Stubbornness

In this subsection we present an interesting class of BBE that exist in all games. In this class, one of the players is “strategically stubborn” in the sense that he plays his undominated Stackelberg strategy (defined below) and has blind beliefs, while his opponent is “flexible” in the sense of having unbiased beliefs.

A strategy is undominated Stackelberg if it maximizes a player's payoff in a setup in which the player can commit to an undominated strategy, and his opponent reacts by choosing the best reply that maximizes player  $i$ 's payoff. Formally:

**Definition 11.** The strategy  $s_i$  is an undominated Stackelberg strategy if it satisfies

$$s_i = \operatorname{argmax}_{s_i \in S_i^U} \left( \max_{s_j \in BR(s_i)} (\pi_i(s_i, s_j)) \right).$$

Let  $\pi_i^{\text{Stac}} = \max_{s_i \in S_i^U} \left( \max_{s_j \in BR(s_i)} (\pi_i(s_i, s_j)) \right)$  be the undominated Stackelberg payoff. Observe that  $\pi_i^{\text{Stac}} \geq \pi_i(s_1^*, s_2^*)$  for any Nash equilibrium  $(s_1^*, s_2^*) \in NE(G)$ .

Our next result shows that every game admits a BBE in which one of the players: (1) has a blind belief, (2) plays his undominated Stackelberg strategy, and (3) obtains his undominated Stackelberg payoff. The opponent has undistorted beliefs. Moreover, this BBE is strong if the undominated Stackelberg strategy is a unique best reply to some undominated strategy of the opponent.

The intuition behind Proposition 7 is as follows. The “strategically stubborn” player  $i$  cannot gain from a deviation, because player  $i$  already obtains the highest possible payoff under the constraint that player  $j$  best-responds to player  $i$ ’s strategy. The “flexible” player  $j$  cannot gain from a deviation, because the “blindness” of player  $i$  implies that player  $i$ ’s behavior remains the same regardless of player  $i$ ’s deviation, and, thus, player  $i$  cannot do better than best-replying to player  $i$ ’s strategy.

**Proposition 7.** *Game  $G = (S, \pi)$  admits a BBE  $((\psi_i^*, Id), (s_i^*, s_j^*))$  for each player  $i$  with the following properties: (1)  $\psi_i^*$  is blind, (2)  $s_i^*$  is an undominated Stackelberg strategy, and (3)  $s_j^* = \max_{s_j \in BR(s_i^*)} (\pi_i(s_i^*, s_j))$ .*

*Moreover,  $((\psi_i^*, Id), (s_i^*, s_j^*))$  is a strong BBE if  $\{s_i^*\} = BR^{-1}(s_j^*)$ .*

*Proof.* Let  $s_i^*$  be an undominated Stackelberg strategy of player  $i$ . Let

$$s_j^* = \operatorname{argmax}_{s_j \in BR(s_i^*)} (\pi_i(s_i^*, s_j)).$$

Let  $\hat{s}_j \in BR^{-1}(s_i^*)$  ( $\{\hat{s}_j\} = BR^{-1}(s_i^*)$  with the additional assumption of the “moreover” part). We now show that  $((\psi_i^* \equiv \hat{s}_j, Id), (s_i^*, s_j^*))$  is a (strong) BBE. It is immediate that  $(s_i^*, s_j^*) \in NE(G_{(\hat{s}_j, Id)})$ , and that both biased beliefs are monotone.

Next, observe that for any biased belief  $\psi_j'$  there is a plausible equilibrium (in any equilibrium) of the biased game  $G_{(\hat{s}_j, \psi_j')}$  in which player  $i$  plays  $s_i^*$ , and player  $j$  gains at most  $\pi_j(s_i^*, s_j^*)$ , which implies that the deviation to  $\psi_j'$  is not profitable to player  $j$  in this plausible equilibrium (in any equilibrium) of the new biased game.

If player  $i$  deviates to a biased belief  $\psi_i'$ , then in any equilibrium of the biased game  $G_{(\psi_i', Id)}$  player  $i$  plays some strategy  $s_i'$  and gains a payoff of at most  $\max_{s_j \in BR(s_i')} (\pi_i(s_i', s_j))$ , and this implies that player  $i$ ’s payoff is at most  $\pi_i^{\text{Stac}}$ , and that he cannot gain by deviating. This shows that  $((\hat{s}_j, Id), (s_i^*, s_j^*))$  is a (strong) BBE.  $\square$

We demonstrate this class of equilibria in a Cournot competition.

**Example 5** (*Well-behaved BBE that yields the Stackelberg outcome in a Cournot competition*). Consider the symmetric Cournot game with linear demand in Example 1:  $G = (S, \pi)$ :  $S_i = \mathbb{R}^+$  and  $\pi_i(s_i, s_j) = s_i \cdot (1 - s_i - s_j)$  for each player  $i$ . Then  $((0, Id), (\frac{1}{2}, \frac{1}{4}))$  is a strong well-behaved BBE that induces the Stackelberg outcome  $(\frac{1}{2}, \frac{1}{4})$ , and yields the Stackelberg-leader payoff of  $\frac{1}{8}$  to player 1 and yields the follower payoff of  $\frac{1}{16}$  to player 2. This is because: (1)  $(\frac{1}{2}, \frac{1}{4}) \in NE(G_{(0, Id)})$ , (2) for any biased belief  $\psi_2'$ , player 1 keeps playing  $\frac{1}{2}$  and as a result player 2’s payoff is at most  $\frac{1}{16}$ , and (3) for any biased belief  $\psi_1'$ , player 2 will best-reply to player’s 1 strategy, and thus player 1’s payoff will be at most his Stackelberg payoff of  $\frac{1}{8}$ .

## 7.2 Folk Theorem Results

In this subsection we present various folk theorem results (i.e., general feasibility results) that show that relaxing either of the two requirements in the definition of a BBE (namely, monotonicity and ruling out implausible equilibria) yields little predictive power in various classes of games. Specifically, we show that in those games a strategy profile is a monotone weak BBE outcome (resp., non-monotone strong BBE outcome) essentially iff it is (1) undominated, and (2) induces a payoff above the undominated minmax payoff.

### 7.2.1 Preliminary Definitions

We begin by defining the notions of monotone weak BBE, and of non-monotone strong BBE.

**Definition 12.** A weak BBE  $(\psi^*, s^*)$  is a *monotone weak BBE* if each biased belief  $\psi_i^*$  is monotone for each player  $i$ .

A weak BBE  $(\psi^*, s^*)$  is a *non-monotone strong BBE* if the inequality  $\pi_i(s'_i, s'_j) \leq \pi_i(s_i^*, s_j^*)$  holds for every player  $i$ , every biased belief  $\psi'_i$ , and every strategy profile  $(s'_i, s'_j) \in NE(G_{(\psi'_i, \psi_j^*)})$ .

Note that (1) a monotone weak BBE is a weakening of the notion of a BBE, which relaxes the requirement of ruling out implausible equilibria, and (2) a non-monotone strong BBE is a weakening of the notion of strong BBE, which relaxes the requirement of monotonicity.

### 7.2.2 Folk Theorem Result: Monotone Weak BBE in Finite Games

We say that a finite game  $G$  admits *best replies with full undominated support*, if, for each player  $i$ , there exists an undominated strategy  $s_i \in S_i^U$  with a support that includes all undominated actions, i.e.,  $\text{supp}(s_i) = A_i \cap S_i^U$ . Two classes of games that admit best replies with full undominated support are:

1. *All two-action games.* The reason for this is as follows. If player  $i$  has a dominant action, then, trivially, the dominant action  $a_i$  is an undominated strategy with a support that includes all undominated actions. If player  $i$  does not have a dominant action, then there must be a strategy of the opponent for which the player is indifferent between his two actions, which implies that there exists an undominated strategy with full support.
2. Any game with a totally mixed equilibrium (e.g., a rock-paper-scissors game).

Our next result focuses on finite games that admit best replies with full undominated support, and shows that in such games a strategy profile  $s^*$  is a monotone weak BBE outcome iff (I)  $s^*$  is undominated, and (II) the payoff of  $s^*$  is above the undominated minmax payoff.

The sketch of the proof is as follows. Each player has a blind belief that his opponent plays her part of the Nash equilibrium with full undominated support. This implies that each player is always indifferent between all undominated actions and, as such, can (1) play  $s_i^*$  on the equilibrium path, and (2) play a punishing strategy that guarantees the opponent a payoff of at most her undominated minmax payoff following any deviation of the opponent.

**Proposition 8** (*Folk Theorem result for monotone weak BBE outcomes*). Let  $G$  be a finite game that admits best replies with full undominated support. Then the following two statements are equivalent:

1. Strategy profile  $(s_1^*, s_2^*)$  is a monotone weak BBE outcome.
2. Strategy profile  $(s_1^*, s_2^*)$  is (I) undominated and (II)  $\pi_i(s_1^*, s_2^*) \geq M_i^U$ .

*Proof.* Proposition 3 implies that “1. $\Rightarrow$ 2.” We now show that “2. $\Rightarrow$ 1.” Assume that  $(s_1^*, s_2^*)$  is undominated, and  $\pi_i(s_1^*, s_2^*) \geq M_i^U$ . For each player  $j$ , let  $s_j^p$  be an undominated strategy that guarantees that player  $i$  obtains, at most, his minmax payoff  $M_i^U$ , i.e.,  $s_j^p = \operatorname{argmin}_{s_j \in S_j^U} (\max_{s_i \in S_i} \pi_i(s_i, s_j))$ . For each player  $j$ , let  $s_j^e \in S_j^U$  be a best-reply strategy with full undominated support, i.e.,  $\operatorname{supp}(s_j^e) = A_j \cap S_j^U$ . For each player  $i$ , let  $s_i^d \in BR^{-1}(s_j^e)$ . The fact that  $s_j^e \in BR(s_i^d)$  implies that  $s_j^*, s_j^p \in \Delta(S_j^U) = \Delta(\operatorname{supp}(s_j^e)) \subseteq BR(s_i^d)$ .

We conclude by showing that  $((s_1^d, s_2^d), (s_1^*, s_2^*))$  is a monotone weak BBE (in which both players have blind beliefs). It is immediate that  $(s_1^*, s_2^*) \in NE(s_1^d, s_2^d)$ . Next, observe that for any deviation of player  $i$  to a different biased belief  $\psi'_i$ , there is a Nash equilibrium of the biased game  $G_{(\psi'_i, s_j^e)}$  in which player  $j$  plays  $s_j^p$ , and, as a result, player  $i$  obtains a payoff of at most  $M_i^U$ , which implies that the deviation is not profitable. Thus,  $(s_1^*, s_2^*)$  is a BBE outcome.  $\square$

Proposition 8 suggests that the notion of monotone weak BBE is too weak. The folk theorem result relies on the incumbents “discriminating” against deviators who have exactly the same perceived behavior as the rest of the population: the incumbents of population  $j$  “punish” deviators by playing  $s_j^p$  against them, while continuing to play  $s_j^*$  against the incumbents, even though both the deviators and the incumbents are perceived to behave the same (i.e.,  $\psi_j^*(s_i^e) = \psi_j^*(s_i^*)$ ).

Example 8 in Appendix A.3 demonstrates that the folk theorem result does not necessarily hold for games that do not admit best replies with full undominated support.

### 7.2.3 Folk Theorem Result: Non-Monotone Strong BBE in Interval Games

In this section we show a folk theorem result for strong BBE in a broad family of interval games in which each payoff function  $\pi_i(s_i, s_j)$  is (1) strictly concave in  $s_i$  and (2) weakly convex in  $s_j$ . Examples of such games include Cournot competitions, price competitions with differentiated goods, public good games, and Tullock contests.

The following result shows that in this class of interval games, any undominated strategy profile  $(s_1^*, s_2^*)$  that induces each player a payoff strictly above the player’s undominated minmax payoff can be implemented as an outcome of a strong BBE. Formally:

**Proposition 9.** Let  $G = (S, \pi)$  be an interval game. Assume that for each player  $i$ ,  $\pi_i(s_i, s_j)$  is strictly concave in  $s_i$  and weakly convex in  $s_j$ . If  $(s_1^*, s_2^*)$  is undominated and  $\pi_i(s_1^*, s_2^*) > M_i^U$  for each player  $i$ , then  $(s_1^*, s_2^*)$  is a non-monotone strong BBE outcome.

The sketch of the proof is as follows (the formal proof is presented in Appendix F.1).

Each player  $j$  has a biased belief  $\psi_j^*$  that (I) distorts  $s_i^*$  into  $BR^{-1}(s_j^*)$ , and (II) distorts any  $s'_i$  that is not in a small neighborhood of  $s_i^*$ , to  $BR^{-1}(s_j^p)$ , where  $s_j^p$  is a “punishing” strategy that



guarantees that player  $i$  obtains at most his undominated minmax payoff. Part (I) implies that  $(s_1^*, s_2^*)$  is an equilibrium of the biased game. Part (II) implies that following any deviation of player  $i$  to a different biased belief, if player  $i$  plays a strategy that is not in a small neighborhood of  $s_i^*$ , then player  $i$  loses from the deviation. Finally, the assumption that the payoff function  $\pi_i(s_i, s_j)$  is convex in  $s_j$  implies that we can “complete” a continuous description of  $\psi_j^*$  for  $s'_i$  that are in a small neighborhood around  $s_i^*$ , such that a player cannot gain from deviating to playing strategies in this small neighborhood.

#### 7.2.4 Discussion of the Folk Theorem Results

The results of this section show that the notion of weak BBE has little predictive power in the sense that, essentially, any undominated strategy profile with a payoff above the undominated minmax payoff is a weak BBE outcome. Moreover, we show that this multiplicity of BBE outcomes holds in large classes of games also when applying a refinement of monotonicity (Prop. 8), or when applying a refinement of strongness (Prop. 9). By contrast, in Section 6 we show that the combination of two plausible requirements, namely, monotonicity and ruling out implausible equilibria, allows us to achieve sharp predictions for the set of BBE outcomes in various interesting classes of games and for the set of biased beliefs that support these outcomes.

Our folk theorem results have similar properties to the famous folk theorem results for repeated games and sufficiently discounted players (see, e.g., [Fudenberg and Maskin, 1986](#)). This is so because it allows for implicit punishments similar to those used in repeated games in order to sustain equilibria. This is because our model assumes that when a player deviates to a different biased belief his opponent can react to the deviation and deter against it.

Observe that our result has somewhat stronger predictive power than the folk theorem result for repeated games, in the sense that the set of monotone weak BBE in one-shot finite games and the set of non-monotone strong BBE in one-shot interval games are each smaller than the set of subgame-perfect equilibria of repeated games between patient players. In particular, the following strategy profiles can be supported as the subgame-perfect equilibrium outcomes of a repeated game between patient players, but they cannot be the outcome of a weak BBE outcome of a one-shot game: (1) strategy profiles in which one of the players plays a strategy that is strictly undominated in the underlying (one-shot) game, and (2) strategy profiles in which some of the players obtain a payoff between the standard minmax payoff and the (higher) undominated minmax payoff.

In Appendix D we show that if one relaxes the assumption that the biased beliefs must be continuous, then one can obtain a folk theorem result in broader classes of games, namely, (1) in all finite games, and (2) in all interval games with strictly concave payoffs.

## 8 Conclusion

Decision makers’ preferences and beliefs may intermingle. In strategic environments distorted beliefs can take the form of a self-serving commitment device. Our paper introduces a formal model for the persistence of such beliefs and proposes an equilibrium concept that supports them. Our analysis characterizes BBE in a variety of strategic environments, such as games with strategic complements

and games with strategic substitutes. In particular, we show that agents present wishful thinking in all BBE in both of these common environments.

Our analysis here deals with simultaneous games of complete information, but the idea of strategically distorted beliefs may play an important role also in sequential games and in Bayesian games. In these frameworks, belief distortion may violate Bayesian updating, and our concept here can potentially offer a theoretical foundation for some of the cognitive biases relating to belief updating. It can also potentially identify the strategic environments in which these biases are likely to occur. We view this as an important research agenda that we intend to undertake in the future.

A different research track that might shed more light on strategic belief distortion is the experimental one. Laboratory experiments often conduct belief elicitation with the support of incentives for truthful revelation. Strong evidence for strategic belief bias in experimental games can be obtained by showing that players assign different beliefs to the behavior of their own counterpart in the game and to a person playing the same role with someone else. In general, our model predicts that beliefs about a third party's behavior are more aligned with reality than those involving one's counterpart in the game. Laboratory experiments can also test whether specific types of belief distortions (such as wishful thinking) arise in the strategic environments that are predicted by our model.

Finally, we point out that strategic beliefs may play an important role in the design of mechanisms and contracts. Belief distortions may destroy the desirable equilibrium outcomes that a standard mechanism aims to achieve. Mechanisms that either induce unbiased beliefs or adjust the rules of the game to account for possible belief biases are expected to perform better.

## References

- ACEMOGLU, D., AND M. YILDIZ (2001): "Evolution of perceptions and play," mimeo.
- ALGER, I., AND J. W. WEIBULL (2013): "Homo moralis, preference evolution under incomplete information and assortative matching," *Econometrica*, 81(6), 2269–2302.
- ATTANASI, G., AND R. NAGEL (2008): "A survey of psychological games: Theoretical findings and experimental evidence," in *Games, Rationality and Behavior: Essays on Behavioral Game Theory and Experiments*, ed. by A. Innocenti, and P. Sbriglia, pp. 204–232. London: Palgrave Macmillan.
- AUMANN, R., AND A. BRANDENBURGER (1995): "Epistemic conditions for Nash equilibrium," *Econometrica*, pp. 1161–1180.
- BABAD, E., AND Y. KATZ (1991): "Wishful thinking: Against all odds," *Journal of Applied Social Psychology*, 21(23), 1921–1938.
- BABCOCK, L., AND G. LOEWENSTEIN (1997): "Explaining bargaining impasse: The role of self-serving biases," *Journal of Economic Perspectives*, 11(1), 109–126.
- BAR-HILLEL, M., AND D. BUDESCU (1995): "The elusive wishful thinking effect," *Thinking & Reasoning*, 1(1), 71–103.

- BARBER, B. M., AND T. ODEAN (2001): “Boys will be boys: Gender, overconfidence, and common stock investment,” *The Quarterly Journal of Economics*, 116(1), 261–292.
- BASU, K. (1994): “The traveler’s dilemma: Paradoxes of rationality in game theory,” *American Economic Review*, 84(2), 391–395.
- BATTIGALLI, P., AND M. DUFWENBERG (2007): “Guilt in games,” *American Economic Review*, 97(2), 170–176.
- (2009): “Dynamic psychological games,” *Journal of Economic Theory*, 144(1), 1–35.
- BATTIGALLI, P., M. DUFWENBERG, AND A. SMITH (2015): “Frustration and anger in games,” mimeo.
- BATTIGALLI, P., AND D. GUAITOLI (1997): “Conjectural equilibria and rationalizability in a game with incomplete information,” in *Decisions, Games and Markets*, pp. 97–124. Berlin: Springer.
- BOLTON, G. E., AND A. OCKENFELS (2000): “ERC: A theory of equity, reciprocity, and competition,” *American Economic Review*, 90(1), 166–193.
- BUDESCU, D. V., AND M. BRUDERMAN (1995): “The relationship between the illusion of control and the desirability bias,” *Journal of Behavioral Decision Making*, 8(2), 109–125.
- BULOW, J. I., J. D. GEANAKOPLOS, AND P. D. KLEMPERER (1985): “Multimarket oligopoly: Strategic substitutes and complements,” *Journal of Political Economy*, 93(3), 488–511.
- CAMERER, C. F., T.-H. HO, AND J.-K. CHONG (2004): “A cognitive hierarchy model of games,” *The Quarterly Journal of Economics*, 119(3), 861–898.
- COSTA-GOMES, M., V. P. CRAWFORD, AND B. BROSETA (2001): “Cognition and behavior in normal-form games: An experimental study,” *Econometrica*, 69(5), 1193–1235.
- DEKEL, E., J. C. ELY, AND O. YILANKAYA (2007): “Evolution of preferences,” *Review of Economic Studies*, 74(3), 685–704.
- DIXIT, A. (1987): “Strategic behavior in contests,” *The American Economic Review*, 77(5), 891–898.
- DOBSON, K., AND R.-L. FRANCHE (1989): “A conceptual and empirical review of the depressive realism hypothesis,” *Canadian Journal of Behavioural Science*, 21(4), 419–433.
- DUFWENBERG, M., AND W. GÜTH (1999): “Indirect Evolution vs. Strategic Delegation: A Comparison of Two Approaches to Explaining Economic Institutions,” *European Journal of Political Economy*, 15(2), 281–295.
- ESPONDA, I. (2013): “Rationalizable conjectural equilibrium: A framework for robust predictions,” *Theoretical Economics*, 8(2), 467–501.
- ESPONDA, I., AND D. POUZO (2016): “Berk–Nash equilibrium: A framework for modeling agents with misspecified models,” *Econometrica*, 84(3), 1093–1130.

- EYSTER, E., AND M. RABIN (2005): “Cursed equilibrium,” *Econometrica*, 73(5), 1623–1672.
- FEHR, E., AND K. M. SCHMIDT (1999): “A theory of fairness, competition, and cooperation,” *Quarterly Journal of Economics*, 114(3), 817–868.
- FERSHTMAN, C., AND U. GNEEZY (2001): “Strategic delegation: An experiment,” *RAND Journal of Economics*, 32(2), 352–368.
- FERSHTMAN, C., K. L. JUDD, AND E. KALAI (1991): “Observable contracts: Strategic delegation and cooperation,” *International Economic Review*, 32(3), 551–559.
- FERSHTMAN, C., AND Y. WEISS (1998): “Social rewards, externalities and stable preferences,” *Journal of Public Economics*, 70(1), 53–73.
- FORBES, D. P. (2005): “Are some entrepreneurs more overconfident than others?,” *Journal of Business Venturing*, 20(5), 623–640.
- FRIEDMAN, D., AND N. SINGH (2009): “Equilibrium vengeance,” *Games and Economic Behavior*, 66(2), 813–829.
- FUDENBERG, D., AND D. K. LEVINE (1993): “Self-confirming equilibrium,” *Econometrica*, 61(3), 523–545.
- FUDENBERG, D., AND E. MASKIN (1986): “The folk theorem in repeated games with discounting or with incomplete information,” *Econometrica*, 54(3), 533–554.
- GANNON, K., AND H. ZHANG (2017): “Evolutionary Justifications for Overconfidence,” mimeo.
- GEANAKOPOLOS, J., D. PEARCE, AND E. STACCHETTI (1989): “Psychological games and sequential rationality,” *Games and Economic Behavior*, 1(1), 60–79.
- GÜTH, W. (1995): “An evolutionary approach to explaining cooperative behavior by reciprocal incentives,” *International Journal of Game Theory*, 24(4), 323–344.
- GÜTH, W., AND S. NAPEL (2006): “Inequality aversion in a variety of games: An indirect evolutionary analysis,” *The Economic Journal*, 116, 1037–1056.
- GÜTH, W., AND M. YAARI (1992): “Explaining reciprocal behavior in simple strategic games: An evolutionary approach,” in *Explaining Process and Change: Approaches to Evolutionary Economics*, ed. by U. Witt, pp. 23–34. Ann Arbor: University of Michigan Press.
- GUTTMAN, J. M. (2003): “Repeated interaction and the evolution of preferences for reciprocity,” *The Economic Journal*, 113(489), 631–656.
- HEIFETZ, A., E. SEGEV, ET AL. (2004): “The evolutionary role of toughness in bargaining,” *Games and Economic Behavior*, 49(1), 117–134.
- HEIFETZ, A., C. SHANNON, AND Y. SPIEGEL (2007a): “The dynamic evolution of preferences,” *Economic Theory*, 32(2), 251–286.

- (2007b): “What to maximize if you must,” *Journal of Economic Theory*, 133(1), 31–57.
- HELLER, Y. (2014): “Overconfidence and diversification,” *American Economic Journal: Microeconomics*, 6(1), 134–153.
- HELLER, Y., AND E. MOHLIN (2017): “Coevolution of deception and preferences: Darwin and Nash meet Machiavelli,” mimeo.
- HELLER, Y., AND D. STURROCK (2017): “Commitments and partnerships,” mimeo.
- HELLER, Y., AND E. WINTER (2016): “Rule rationality,” *International Economic Review*, 57(3), 997–1026.
- HEROLD, F., AND C. KUZMICS (2009): “Evolutionary stability of discrimination under observability,” *Games and Economic Behavior*, 67, 542–551.
- HOLMSTROM, B. (1982): “Moral hazard in teams,” *The Bell Journal of Economics*, 13(2), 324–340.
- INOUE, Y., Y. TONOOKA, K. YAMADA, AND S. KANBA (2004): “Deficiency of theory of mind in patients with remitted mood disorder,” *Journal of Affective Disorders*, 82(3), 403–409.
- JEHIEL, P. (2005): “Analogy-based expectation equilibrium,” *Journal of Economic Theory*, 123(2), 81–104.
- KOÇKESEN, L., E. A. OK, AND R. SETHI (2000): “Evolution of interdependent preferences in aggregative games,” *Games and Economic Behavior*, 31(2), 303–310.
- LORD, C. G., L. ROSS, AND M. R. LEPPER (1979): “Biased assimilation and attitude polarization: The effects of prior theories on subsequently considered evidence,” *Journal of Personality and Social Psychology*, 37(11), 2098–2109.
- MALMENDIER, U., AND G. TATE (2005): “CEO overconfidence and corporate investment,” *Journal of Finance*, 60(6), 2661–2700.
- MAS-COLELL, A., M. D. WHINSTON, AND J. R. GREEN (1995): *Microeconomic Theory*, vol. 1. New York: Oxford university press.
- MAYNARD-SMITH, J., AND G. PRICE (1973): “The logic of animal conflict,” *Nature*, 246, 15–18.
- MAYRAZ, G. (2013): “Wishful thinking,” Discussion paper, University of Melbourne.
- MILGROM, P., AND J. ROBERTS (1990): “Rationalizability, learning, and equilibrium in games with strategic complementarities,” *Econometrica*, 58(6), 1255–1277.
- NAGEL, R. (1995): “Unraveling in guessing games: An experimental study,” *American Economic Review*, 85(5), 1313–1326.
- NYARKO, Y., AND A. SCHOTTER (2002): “An experimental study of belief learning using elicited beliefs,” *Econometrica*, 70(3), 971–1005.

- PALFREY, T. R., AND S. W. WANG (2009): “On eliciting beliefs in strategic games,” *Journal of Economic Behavior & Organization*, 71(2), 98–109.
- POTTERS, J., AND S. SUETENS (2009): “Cooperation in experimental games of strategic complements and substitutes,” *The Review of Economic Studies*, 76(3), 1125–1147.
- RABIN, M. (1993): “Incorporating fairness into game theory and economics,” *American Economic Review*, 83(5), 1281–1302.
- ROSS, L., AND C. ANDERSON (1982): “Shortcomings in attribution processes: On the origins and maintenance of erroneous social judgments,” in *Judgement under Uncertainty: Heuristics and Biases*, ed. by D. Kahnemann, P. Slovic, and A. Tversky, pp. 129–152. Cambridge: Cambridge University Press.
- RUBINSTEIN, A., AND A. WOLINSKY (1994): “Rationalizable conjectural equilibrium: Between Nash and rationalizability,” *Games and Economic Behavior*, 6(2), 299–311.
- RUSTICHINI, A., AND M. C. VILLEVAL (2014): “Moral hypocrisy, power and social preferences,” *Journal of Economic Behavior & Organization*, 107, 10–24.
- SINGH, N., AND X. VIVES (1984): “Price and quantity competition in a differentiated duopoly,” *The RAND Journal of Economics*, 15(4), 546–554.
- STAHL, D. O., AND P. W. WILSON (1994): “Experimental evidence on players’ models of other players,” *Journal of Economic Behavior and Organization*, 25(3), 309–327.
- WINTER, E., I. GARCIA-JURADO, AND L. MENDEZ-NAYA (2017): “Mental equilibrium and rational emotions,” *Management Science*, 63(5), 1302–1317.

## Online Appendices

### A Additional Examples

#### A.1 A Non-Nash Strong BBE Outcome in a Zero-Sum Game

The following example shows that although the weak BBE payoff must be the Nash equilibrium payoff in a zero-sum game, the strategy profile sustaining it need not be a Nash equilibrium.

**Example 6.** Consider the symmetric rock–paper–scissors zero-sum game described in Table 2. We show that  $\left(\left(I_d, \left(\frac{1}{3}, \frac{1}{3}, \frac{1}{3}\right)\right), \left(R, \left(\frac{1}{3}, \frac{1}{3}, \frac{1}{3}\right)\right)\right)$  is a strong BBE, in which the player 1 (he) has

Table 2: Symmetric Rock-Paper-Scissors Zero-Sum Game Payoffs

	R	P	S
R	0, 0	0, 1	1, 0
P	1, 0	0, 0	0, 1
S	0, 1	1, 0	0, 0

undistorted beliefs and plays  $R$ , while player 2 (she) has a blind belief that the opponent always mixes equally, and she mixes equally. It is immediate that  $\left(R, \left(\frac{1}{3}, \frac{1}{3}, \frac{1}{3}\right)\right) \in NE\left(G_{\left(I_d, \left(\frac{1}{3}, \frac{1}{3}, \frac{1}{3}\right)\right)}\right)$ , and the equilibrium payoff to each player is zero. Next, observe that after any deviation of player 1 to a biased belief  $\psi'_1$ , there is an equilibrium of the game  $G_{\left(\psi'_1, \left(\frac{1}{3}, \frac{1}{3}, \frac{1}{3}\right)\right)}$  in which player 2 mixes equally and player 1 obtains a payoff of zero. Finally, observe that after any deviation of player 2 to a biased belief  $\psi'_2$ , player 1 obtains a payoff of at least zero (her minmax payoff in  $G_{\left(I_d, \psi'_2\right)}$ ) in any Nash equilibrium in  $G_{\left(I_d, \psi'_2\right)}$ , which implies that player 2 obtains a payoff of at most zero, and, as a result, she does not gain from the deviation.

#### A.2 Prisoner’s Dilemma with a Weakly Dominated Withdrawal Strategy

Proposition 2 implies, in particular, that defection is the unique weak BBE outcome in the prisoner’s dilemma game. The following example demonstrates that a relatively small change to the prisoner’s dilemma game, namely, adding a third weakly dominated “withdrawal” strategy that transforms “cooperation” into a weakly dominated strategy, can allow us to sustain cooperation as a strong BBE outcome. This is done by means of biases under which a player believes that his opponent is planning to withdraw from the game whenever he intends to cooperate, which makes cooperation a rational move.

**Example 7.** Consider the variant of the prisoner’s dilemma game with a third “withdrawal” action as described in Table 3. In this symmetric game both players get a high payoff of 10 if they both play action  $c$  (interpreted as cooperation). If one player plays  $d$  (*defection*) and his opponent plays  $c$ , then the defector gets 11 and the cooperator gets 0. If both players defect, then each of them gets a payoff of 1. Finally, if either player plays action  $w$  (interpreted as *withdrawal*), then both players get 0. Observe that defection is a weakly dominant action, and that the game admits two Nash equilibria:  $(w, w)$  and  $(d, d)$ , inducing respective symmetric payoffs of zero and one.



Table 3: Prisoner's Dilemma Game with a Withdrawal Action

	$c$	$d$	$w$
$c$	10,10	0,11	0,0
$d$	11,0	1,1	0,0
$w$	0,0	0,0	0,0

We identify a mixed action with a vector  $(\alpha_c, \alpha_d, \alpha_w)$ , where  $\alpha_c \geq 0$  (resp.,  $\alpha_d \geq 0$ ,  $\alpha_w \geq 0$ ) denotes the probability of choosing action  $c$  (resp.,  $d$ ,  $w$ ). For each player  $i$ , let  $\psi_i$  be the following biased-belief function:

$$\psi_i^*(\alpha_c, \alpha_d, \alpha_w) = (0, \alpha_d, \alpha_c + \alpha_w).$$

We now show that  $((\psi_1^*, \psi_2^*), (c, c))$  is a non-monotone strong BBE in which both players obtain a high payoff of 10 (which is strictly better than the best Nash equilibrium payoff, and strictly better than the Stackelberg payoff of each player). Observe first that  $c \in BR(\psi_i^*(c)) = BR(w)$ , which implies that  $(c, c) \in NE(G_{(\psi_1^*, \psi_2^*)})$ . Next, consider a deviation of player  $i$  to biased belief  $\psi'_i$ . Observe that player  $i$  can gain a payoff higher than 10 only if he plays action  $d$  with positive probability, but this implies that the unique best reply of player  $j$  to his biased belief about player  $i$ 's strategy is defection, which implies that player  $i$  obtains a payoff of at most one.

### A.3 The Folk Theorem Result Does not Hold for All Finite Games

The following example demonstrates that the folk theorem result (Proposition 4) does not necessarily hold for games that do not admit best replies with full undominated support.

**Example 8.** Consider the three-action symmetric game described in Table 4. Observe that all

Table 4: A Game in which  $(a, a)$  is not a Monotone Weak BBE Outcome

	$a$	$b$	$c$
$a$	2, 2	2, 3	1.1, 3
$b$	3, 2	3, 3	1, 0
$c$	3.1, 1	0, 1	0, 0

the actions in the game are undominated, and that the game does not admit best replies with full undominated support: there is no strategy of the opponent for which one of the players has a best reply with full support. This is so because action  $a$  ( $c$ ) is a best reply only to his opponent's strategies that assign a probability of at least 90% to action  $c$  ( $a$ ), which implies that actions  $a$  and  $c$  cannot be best replies simultaneously. Observe that the undominated minmax payoff of each player is equal to 1 (because the opponent can play the undominated action  $c$ , and by playing this the opponent guarantees that the player gets a payoff of at most 1).

Consider the undominated action profile  $(a, a)$  (which induces a payoff strictly above the undominated minmax payoff to each player). We will show that  $(a, a)$  is not a monotone weak BBE (which demonstrates that the folk theorem result of Proposition 8 does not hold in this game). Assume to the contrary that  $(a, a)$  is a monotone weak BBE. Let  $((\psi_1^*, \psi_2^*), (a, a))$  be a monotone weak BBE.

The fact that  $(a, a) \in NE(G_{(\psi_1^*, \psi_2^*)})$  implies that  $\psi_1^*(a)(c) > 90\%$ . Consider a deviation of player 2 to having the blind belief  $\psi_2' = b$ . Observe that player 2 plays action  $b$  in any equilibrium of  $G_{(\psi_1^*, \psi_2')}$ . The monotonicity of  $\psi_1^*$  implies that  $\psi_1^*(b)(a) \leq \psi_1^*(a)(a) \leq 1 - \psi_1^*(a)(c) \leq 10\%$ , which implies that the best reply of player 1 to the perceived strategy of player 2 ( $\psi_1^*(b)$ ) does not have action  $c$  in its support. This implies that player 1 gains a payoff of at least 3 in any Nash equilibrium of the new biased game  $G_{(\psi_1^*, \psi_2')}$ , which contradicts  $((\psi_1^*, \psi_2^*), (a, a))$  being a monotone weak BBE.

#### A.4 Examples of Games with Strategic Complements

In this subsection we analyze three examples of games with strategic complements: input games, stag hunt games, and the traveler's dilemma.

Our first example demonstrates how to implement the undominated Pareto optimal profile as a strong BBE in an input (or partnership game).

**Example 9** (*Input games*). Consider the following input game (closely related games are analyzed in, among others, [Holmstrom, 1982](#) and [Heller and Sturrock, 2017](#)). Let  $S_i = S_j = [0, 1]$ , and let the payoff function be  $\pi_i(s_i, s_j, \rho) = s_i \cdot s_j - \frac{s_i^2}{2\rho}$ , where the parameter  $\frac{1}{\rho}$  is interpreted as the cost of effort. One can show that (1) the best-reply function of each agent is to exert an effort that is  $\rho < 1$  times smaller than the opponent's (i.e.,  $BR(s_j) = \rho \cdot s_j$ ), (2) in the unique Nash equilibrium each player exerts no effort  $s_i = s_j = 0$ , (3) the highest undominated strategy of each player  $i$  is  $s_i = \rho$ , and (4) the undominated strategy profile  $(\rho, \rho)$  is Nash improving and yields the best payoff to both players out of all the undominated symmetric strategy profiles. Let  $\psi_i^*$  be the following biased-belief function:

$$\psi_i^*(s_j) = \begin{cases} \frac{s_j}{\rho} & s_j < \rho \\ 1 & s_j \geq \rho. \end{cases}$$

Observe that  $\psi_i^*$  is monotone and exhibits wishful thinking. We now show that  $((\psi_1^*, \psi_2^*), (\rho, \rho))$  is a strong BBE. Observe that  $BR(\psi_i^*(s_j)) = BR(\frac{s_j}{\rho}) = s_j$  for any  $s_j \leq \rho$ , and that  $BR(\psi_i^*(s_j)) = BR(1) = \rho$  for any  $s_j \geq \rho$ . This implies that  $(\rho, \rho) \in NE(G_{(\psi_1^*, \psi_2^*)})$ , and that for any player  $i$ , any biased belief  $\psi_i'$ , and any Nash equilibrium  $(s'_1, s'_2)$  of the biased game  $G_{(\psi_i', \psi_j)}$ ,  $s'_j = \min(s'_i, \rho)$ . This implies that  $\pi_i(s'_1, s'_2) \leq \pi_i(\rho, \rho)$ , which shows that  $((\psi_1^*, \psi_2^*), (\rho, \rho))$  is a strong BBE. Observe that this BBE induces only a small distortion in the belief of each player, assuming that  $\rho$  is sufficiently close to one:

$$|\psi_i^*(s_j) - s_j| < \left| \frac{s_j}{\rho} - s_j \right| < \frac{1 - \rho}{\rho}.$$

Our second example characterizes the set of BBE outcomes (and their supporting beliefs) in stag hunt games.

**Example 10** (*Stag hunt games*). Stag hunt is a two-action game describing a conflict between safety and social cooperation. Specifically, each player  $i$  has two actions:  $s_i$  ("stag") and  $h_i$  ("hare"), and his ordinal preferences are  $(s_i, s_j) \succ_i (h_i, s_j) \succeq_i (h_i, h_j) \succ_i (s_i, h_j)$ . Table 5 presents the payoff of a general stag hunt game, where we have normalized, without loss of generality, the payoff of each player when playing action profile  $(s_i, s_j)$  ( $(h_i, h_j)$ ) to be one (zero), and where each  $g_i$  is positive and each  $l_i$  is in the interval  $(0, 1)$ . A common interpretation of stag hunt games (à la Jean-Jacques

Table 5: Stag Hunt Game ( $g_1, g_2 \in (0, 1]$  and  $l_1, l_2 > 0$ )

	$s_2$	$h_2$
$s_1$	1, 1	$-l_1, g_2$
$h_1$	$g_1, -l_1$	0, 0

Rousseau) is a situation in which two individuals go hunting. Each can individually choose to hunt a stag or to hunt a hare. Each player must choose an action without knowing the choice of the other. If an individual hunts a stag, he must have the cooperation of his opponent in order to succeed. An individual can get a hare by himself, but a hare is worth less than a stag. It is well known that the game admits 3 equilibria:  $(s_i, s_j)$ ,  $(h_i, h_j)$ , and  $(\alpha_1^*, \alpha_2^*)$ , with

$$\alpha_i^* = \frac{l_j}{l_j + (1 - g_j)} \in (0, 1),$$

where each  $\alpha_i$  represents the probability that player  $i$  plays  $s_i$ .

Applying the analysis of the previous section shows that the game admits 3 classes of BBE:

- Hunting the hare:  $((\psi_1^*, \psi_2^*), (0, 0))$ , where each  $\psi_i^*$  is an arbitrary monotone biased belief that satisfies  $\psi_i^*(1) \geq \alpha_i^*$ .
- Hunting the stag:  $((\psi_1^*, \psi_2^*), (1, 1))$ , where each  $\psi_i^*$  is an arbitrary monotone biased belief that satisfies  $\psi_i^*(1) \leq \alpha_i^*$ .
- Mixing with less weight to hunting the stag, wishful thinking, and responsiveness to bad news:  $((\psi_1^*, \psi_2^*), (\beta_1, \beta_2))$ , where for each player  $i$ : (1) the payoff is above the minmax payoff:  $\pi_i(\beta_i, \beta_j) \geq 0$ , (2) the players hunt the stag less often in the unique Nash equilibrium:  $\beta_i \in (0, \alpha_i^*)$ , (3) wishful thinking:  $\psi_i^*(\beta_j) = \alpha_j^* > \beta_j$ , (4) responsiveness to bad news:  $\psi_i^*(\alpha) = \alpha_j^*$  for each  $\alpha \geq \beta_j$ , and  $\psi_i^*(\alpha) < \alpha_j^*$  for each  $\alpha < \beta_j$ .

Observe that any profile  $(\beta_1, \beta_2)$ , where  $\beta_i \in (\alpha_i^*, 1)$ , cannot be a BBE outcome. If  $\beta_j = 1$ , then player  $i$  can gain by deviating to  $\psi_i' \equiv 1$ , as the unique equilibrium of the new biased game is  $(1, 1)$ , which induces a higher payoff to player  $i$  relative to  $(\beta_i, \beta_j)$ . If  $\beta_j < 1$ , then player  $j$  can gain by deviating to  $\psi_j' \equiv 1$ , as the only possible equilibria of the new biased game are  $(1, 1)$  and  $(\beta_i, 1)$ , both of which induce a higher payoff to player  $j$  relative to  $(\beta_i, \beta_j)$ .

Our third example deals with the traveler's dilemma game, in which each agent has 100 pure ordered actions that have a discrete payoff structure that resembles strategic complementarity in interval games. We demonstrate how to implement the undominated Pareto optimal profile in this game as a strong BBE outcome that presents wishful thinking.

**Example 11** (*Implementing the undominated Pareto optimal profile as a strong BBE in the traveler's dilemma*).

Consider the following version of the traveler's dilemma game (Basu, 1994). Each player has 100 actions ( $A_i = \{1, \dots, 100\}$ ), and the payoff function of each player is

$$\pi_i(a_i, a_j) = \begin{cases} a_i + 2 & a_i < a_j \\ a_i & a_i = a_j \\ a_j - 2 & a_i > a_j. \end{cases}$$

The interpretation of the game is as follows. Two identical suitcases have been lost, each owned by one of the players. Each player has to evaluate the value of his own suitcase. Both players get a payoff equal to the minimal evaluation (as the suitcases are known to have identical values), and, in addition, if the evaluations differ, then the player who gave the lower (higher) evaluation gets a bonus (malus) of 2 to his payoff.

It is well known that the unique Nash equilibrium is (1, 1), which yields a low payoff of one to each player. Observe that the traveler's dilemma has positive spillovers, in the sense that it is always weakly better for a player if his opponent chooses a higher action. The traveler's dilemma has strategic complementarity in the sense that the best reply of an agent is to stop one stage before his opponent, and, thus, an agent has an incentive to choose a higher action if his opponent chooses a higher action.

Observe that action 99 is the “highest” undominated action of each player (as 99 is a best reply against 100, and as action 100 is not a best reply against any of the opponent's strategies). In what follows, we construct a strong BBE exhibiting wishful thinking that yields a payoff of 99 to each player in the undominated symmetric Pareto-optimal strategy profile.

We define the biased belief  $\psi_i^*$  as follows:

$$\psi_i^*(\alpha_1, \alpha_2, \dots, \alpha_{99}, \alpha_{100}) = \left( \alpha_1, \alpha_2, \dots, \frac{\alpha_{99}}{2}, \frac{\alpha_{99}}{2} + \alpha_{100} \right).$$

In what follows we show that  $((\psi_1^*, \psi_2^*), (99, 99))$  is a strong BBE. Observe first that  $\psi_1^*(99) = (0, \dots, 0, \frac{1}{2}, \frac{1}{2})$ , which implies that  $99 \in BR(\psi_i^*(99))$ , and, thus,  $(99, 99) \in NE(G_{(\psi_1^*, \psi_2^*)})$ . Let  $\psi'_1$  be an arbitrary perception bias of player  $i$ . Observe that player  $i$  never plays action 100 in a any Nash equilibrium of any biased game, because action 100 is not a best reply against any strategy of player  $j$ . Next observe that player  $i$  can obtain a payoff higher than 99 only if (1) player  $j$  chooses action 99 with a positive probability, and (2) player  $i$  chooses action 98 with a probability strictly higher than his probability of playing action 100. However, the biased belief  $\psi_j^*$  of player  $j$  implies that if player  $i$  chooses action 98 with a probability strictly higher than his probability of playing 100, then player  $j$  never chooses action 99 in any Nash equilibrium of the induced biased game because action 99 yields a strictly lower payoff to player  $j$  than action 98 against the perceived strategy of player  $i$  (because according to this perceived strategy, player  $i$  plays action 100 with a probability strictly less than player  $i$ 's probability of playing either action 98 or action 99).

Note that the BBE equilibrium outcome (99, 99) is consistent with level-1 behavior in the level- $k$  and cognitive hierarchy literature (see, e.g., Stahl and Wilson, 1994; Nagel, 1995; Costa-Gomes, Crawford, and Broseta, 2001; Camerer, Ho, and Chong, 2004), according to which each agent believes

that his opponent is following a focal non-strategic action (the action 100 in the traveler’s dilemma), and best-replies to this belief. The notion of BBE can help explain why such level-k behavior induces a strategic advantage in the long run, and why, therefore, it is likely to emerge in an equilibrium.

### A.5 Hawk-Dove Game

The following example characterizes the set of BBE (and their supporting beliefs) in a hawk-dove game (which is a game of strategic substitutes).

**Example 12** (*The Hawk-dove game*). The hawk-dove (or “chicken”) game is a two-action game in which each player  $i$  has two actions:  $d_i$  (interpreted as a “dove”-like action of willingness to share a resource with the opponent) and  $h_i$  (interpreted as a “hawk”-like action of insistence on getting the whole resource, even if this requires fighting against the opponent), and where the ordinal preferences of each player  $i$  are  $(h_i, d_j)$  (getting the resource)  $\succ (d_i, d_j)$  (sharing the resource)  $\succ (d_i, h_j)$  (not getting the resource)  $\succ (h_i, h_j)$  (being involved in a serious fight). Table 6 presents the payoff of a general two-action hawk-dove game, where we have normalized, without loss of generality, the payoff of each player when playing action profile  $(d_i, d_j)$  ( $(h_i, h_j)$ ) to be one (zero), and where each  $g_i$  positive and each  $l_i$  is in the interval  $(0, 1)$ .

Table 6: Hawk-Dove Game ( $g_1, g_2 > 0$  and  $l_1, l_2 \in (0, 1)$ )

	$d_2$	$h_2$
$d_1$	1, 1	$1 - l_1, 1 + g_2$
$h_1$	$1 + g_1, 1 - l_1$	0, 0

It is well known that the hawk-dove game admits three equilibria: two pure equilibria  $(d_1, h_2)$  and  $(h_1, d_2)$ , and one mixed equilibrium  $(\alpha_1^*, \alpha_2^*)$ , where the probability that player  $i$  plays action  $\alpha_i^*$  is

$$\alpha_i^* = \frac{1 - l_j}{g_j + (1 - l_j)} \in (0, 1), \quad \text{and} \quad \pi(\alpha_i^*, \alpha_j^*) = \alpha_j^* \cdot (1 + g_i) = 1 - \frac{g_i}{g_i + (1 - l_i)} \cdot l_i.$$

The undominated minmax payoff of each player coincides with the minmax payoff of each player (as there are no dominated actions), and it is equal to  $M_i^U = 1 - l_i$ , which is obtained when the opponent plays  $h_j$ .

Applying the analysis of the previous section shows that the game admits 3 classes of BBE:

- Pure equilibrium hawk-dove:  $((\psi_i^*, \psi_j^*), (0, 1))$ , where (1)  $\psi_i^*$  is an arbitrary monotone biased belief that satisfies  $\psi_i^*(0) \geq \alpha_i^*$ , and (2)  $\psi_j^*$  is an arbitrary monotone biased belief that satisfies  $\psi_j^*(0) \leq \alpha_j^*$ .
- Mixing (with less weight to playing dove), wishful thinking, and one-directional blindness:  $((\psi_1^*, \psi_2^*), (\beta_1, \beta_2))$ , where for each player  $i$ : (1) the payoff is above the minmax payoff:  $\pi_i(\beta_i, \beta_j) \geq 1 - l_i$ , (2)  $\beta_i \in (0, \alpha_i^*)$  (i.e., agents play dove less often in the unique Nash

equilibrium), (3) wishful thinking:  $\psi_i^*(\beta_j) = \alpha_j^* > \beta_j$ , and (4) responsiveness only to good news:  $\psi_i^*(\alpha) = \alpha_j^*$  for each  $\alpha \leq \beta_j$ , and  $\psi_i^*(\alpha) > \alpha_j^*$  for each  $\alpha > \beta_j$ .

Observe that any profile  $(\beta_1, \beta_2)$  where  $\beta_i \in (\alpha_i^*, 1)$  cannot be a BBE outcome. If  $\beta_j = 1$ , then player  $i$  can gain by deviating to  $\psi_i' \equiv 1$ , as the unique equilibrium of the new biased game is  $(0_i, 1_j)$ , which induces a higher payoff to player  $i$  relative to  $(\beta_i, \beta_j)$ . If  $\beta_j < 1$ , then player  $j$  can gain by deviating into  $\psi_j' \equiv 1$ , as the only possible equilibria of the new biased game are  $(1_i, 0_j)$  and  $(\beta_i, 1_j)$ , both of which induce a higher payoff to player  $j$  relative to  $(\beta_i, \beta_j)$ .

## B Evolutionary Interpretation of BBE

In this section we present a formal definition of strong BBE that is exactly analogous to the definition of a stable configuration à la [Dekel, Ely, and Yilankaya \(2007\)](#). This shows that our static solution concept of strong BBE captures evolutionary stability in the same way as the solution concepts used in the literature on “indirect evolution of preferences.” Finally, we illustrate a detailed example of a possible learning dynamic that may result in convergence to strong BBE.

### B.1 Evolutionary Definition of Strong BBE à la [Dekel, Ely, and Yilankaya \(2007\)](#)

In this subsection we present a definition of a strong BBE that is completely analogous to the definition of a stable configuration à la [Dekel, Ely, and Yilankaya \(2007\)](#) (henceforth DEY) for the case of perfect observability of the opponent’s type (i.e.,  $p = 1$  in DEY).

In the adaptation of the notion of stable configuration à la [Dekel, Ely, and Yilankaya \(2007\)](#) to our setup we change two aspects (and only these aspects):

1. We deal with general two-player games played between two different populations, rather than DEY’s setup that deals with symmetric two-player games played within a single population.
2. Each agent in DEY’s model is endowed with a type that determines the agent’s subjective preferences. By contrast, in our setup each agent is endowed with a type that determines the agent’s monotone biased belief.
3. We focus on homogeneous configurations. DEY’s general definitions allow one to deal with heterogeneous configurations (in which different incumbents may have different types). However, their results mainly deal with homogeneous configurations (in which all incumbents have the same type). Therefore, to ease notation, we focus on homogeneous configurations in our adaptation of DEY’s definitions.

After adapting DEY’s definition of a homogeneous configuration (page 689 in DEY) to the three aspects mentioned above, their definition is as follows:

**Definition 13.** A (homogeneous) *configuration* is a pair  $((\psi_1^*, \psi_2^*), (s_1^*, s_2^*))$ , where, for each player  $i$ , function  $\psi_i^*$  is a monotone biased belief of player  $i$  and  $s_i^*$  is a strategy of player  $i$  satisfying  $s_i^* \in BR(\psi_i^*(s_j^*))$ .

It is immediate that any monotone weak BBE is a configuration.

Next, DEY present a notion of a balanced configuration (page 689 in DEY) that is trivially satisfied by any homogeneous configuration.

Consider two continuum populations of mass one that follow a configuration  $((\psi_1^*, \psi_2^*), (s_1^*, s_2^*))$ . Assume that one of these populations (say, population  $i$ ) is invaded by a small group of  $0 < \epsilon < 1$  mutants with a different biased belief  $\psi'_i \neq \psi_i^*$ . DEY assume that (1) such a mutation can destabilize a configuration by resulting in the mutants achieving a higher fitness than the incumbents of the same population<sup>6</sup>  $i$ , and (2) the incumbents continue to play the same behavior among themselves (what DEY calls “focal equilibria”).

Let  $\Psi_i$  be the set of all biased beliefs of player  $i$ . Following DEY (page 690 in DEY) we define  $N_{i,\epsilon}(\psi_i^*, \psi'_i) \in \Delta(\Psi_i)$  to be the set of distributions over biased beliefs in population  $i$  resulting from entry by no more than  $\epsilon$  mutants. Formally,

$$N_{i,\epsilon}(\psi_i^*, \psi'_i) = \{\mu'_i \in \Delta(\Psi_i) \mid \mu'_i = (1 - \epsilon') \cdot \psi_i^* + \epsilon' \cdot \psi'_i, \epsilon' < \epsilon\}.$$

Given a configuration  $((\psi_1^*, \psi_2^*), (s_1^*, s_2^*))$  and a post-entry distribution of biased beliefs in population  $i$   $\tilde{\mu}_i \in N_{i,\epsilon}(\psi_i^*, \psi'_i)$ , a *post-entry focal configuration* is a pair  $((\tilde{\mu}_i, \psi_j^*), (s'_i, s'_j))$ , where (1)  $s'_i \in BR(\psi'_i(s'_j))$  is interpreted as the mutant’s strategy, and (2)  $s'_j \in BR(\psi_j^*(s'_i))$  is interpreted as population  $j$ ’s strategy against the mutants. The incumbents are assumed to play the same pre-entry strategies  $(s_i^*, s_j^*)$  when being matched among themselves. Let  $B(\tilde{\mu}_i)$  denote the set of all post-entry focal configurations.

Following DEY (Definition 3 on page 691 in DEY), we define DEY-stability of a configuration as follows.

**Definition 14.** Configuration  $((\psi_1^*, \psi_2^*), (s_1^*, s_2^*))$  is *DEY-stable* if there exists  $\epsilon > 0$  such that for every player  $i$ , every biased belief  $\psi'_i$ , every post-entry distribution of biased beliefs  $\tilde{\mu}_i \in N_{i,\epsilon}(\psi_i^*, \psi'_i)$ , and every post-entry focal configuration  $((\tilde{\mu}_i, \psi_j^*), (s'_i, s'_j))$ , the mutants are weakly outperformed relative to the incumbents’ payoff (in their own population), i.e.,  $\pi_i(s'_i, s'_j) \leq \pi_i(s_i^*, s_j^*)$ .

## B.2 Equivalence between the Definitions

The following result shows that the definition of a stable configuration coincides with our definition of strong BBE.

**Proposition 10.** A configuration  $((\psi_1^*, \psi_2^*), (s_1^*, s_2^*))$  is *DEY-stable* iff it is a *strong BBE*.

*Proof.* “If” part: Let  $((\psi_1^*, \psi_2^*), (s_1^*, s_2^*))$  be a strong BBE. Let  $\epsilon > 0$ ,  $i \in \{1, 2\}$ , and  $\psi'_i \in \Psi_i$ . Let  $\tilde{\mu}_i \in N_{i,\epsilon}(\psi_i^*, \psi'_i)$  be a post-entry distribution of biased beliefs. Let  $((\tilde{\mu}_i, \psi_j^*), (s'_i, s'_j))$  be a post-entry focal configuration. The fact that  $((\tilde{\mu}_i, \psi_j^*), (s'_i, s'_j))$  is a post-entry focal configuration implies that  $s'_i \in BR(\psi'_i(s'_j))$  and  $s'_j \in BR(\psi_j^*(s'_i))$ . The fact that it is a strong BBE implies that  $\pi_i(s'_i, s'_j) \leq \pi_i(s_i^*, s_j^*)$ , which shows that  $((\psi_1^*, \psi_2^*), (s_1^*, s_2^*))$  is DEY-stable.

---

<sup>6</sup>Under imperfect observability, a mutant can destabilize a configuration by unraveling the original equilibrium behavior, thereby causing the incumbents’ strategies to substantially diverge following the mutant’s entry into the population. This cannot happen under perfect observability, as the incumbents can always exhibit the same equilibrium behavior when being matched against other incumbents (see, page 690 in DEY for a discussion of focal equilibria).



“Only if” part: Let  $((\psi_1^*, \psi_2^*), (s_1^*, s_2^*))$  be DEY-stable configuration. Let  $i \in \{1, 2\}$  and  $\psi'_i \in \Psi_i$ . Let  $(s'_i, s'_j) \in NE\left(G_{(\psi'_i, \psi_j^*)}\right)$  be an equilibrium of the new biased game. Let  $\epsilon > 0$ . Let  $\tilde{\mu}_i \in N_{i,\epsilon}(\psi_i^*, \psi'_i)$  be a post-entry distribution of biased beliefs. For each  $(s'_i, s'_j) \in NE\left(G_{(\psi'_i, \psi_j^*)}\right)$ , let  $((\tilde{\mu}_i, \psi_j^*), (s'_i, s'_j))$  be a post-entry focal configuration. The assumption that  $((\psi_1^*, \psi_2^*), (s_1^*, s_2^*))$  is DEY-stable implies that  $\pi_i(s'_i, s'_j) \leq \pi_i(s_i^*, s_j^*)$ . This implies that  $((\psi_1^*, \psi_2^*), (s_1^*, s_2^*))$  is a strong BBE.  $\square$

*Remark* (Allowing multiple simultaneous invasions of mutants). The definition of DEY-stability presented above is unaffected when various groups of mutants simultaneously invade one of the populations. By contrast, if one were to require a stable configuration to resist simultaneous invasions of two groups of mutants, one invasion of each population, it would require a refinement of the concept of strong BBE, in the spirit of [Maynard-Smith and Price’s \(1973\)](#) notion of evolutionary stability, such that if both  $\psi'_1$  and  $\psi'_2$  are best replies against configuration  $((\psi_1^*, \psi_2^*), (s_1^*, s_2^*))$ , then (1)  $\psi_1^*$  should be a strictly better reply against  $\psi'_2$  (relative to  $\psi'_1$ ), and (2)  $\psi_2^*$  should be a strictly better reply against  $\psi'_1$  (relative to  $\psi'_2$ ).

Similarly, one can formulate a definition of stability equivalent to that of monotone BBE by requiring the mutants to be weakly outperformed in at least one post-entry focal configuration.

**Definition 15.** Configuration  $((\psi_1^*, \psi_2^*), (s_1^*, s_2^*))$  is *weakly stable* if there exists  $\epsilon > 0$  such that for every player  $i$ , every biased belief  $\psi'_i$ , and every post-entry distribution of biased beliefs  $\tilde{\mu}_i \in N_{i,\epsilon}(\psi_i^*, \psi'_i)$ , there exists a post-entry focal configuration  $((\tilde{\mu}_i, \psi_j^*), (s'_i, s'_j))$  in which the mutants are weakly outperformed relative to the incumbents’ payoff, i.e.,  $\pi_i(s'_i, s'_j) \leq \pi_i(s_i^*, s_j^*)$ .

The following result shows that the definition of a weakly stable configuration coincides with our definition of weak BBE. The simple proof, which is analogous to the proof of [10](#), is omitted for brevity.

**Proposition 11.** A configuration  $((\psi_1^*, \psi_2^*), (s_1^*, s_2^*))$  is weakly stable iff it is a monotone weak BBE.

Finally, one can formulate a definition of stability equivalent to that of a BBE by requiring the mutants to be weakly outperformed in at least one plausible post-entry focal configuration.

**Definition 16.** Given configuration  $((\psi_1^*, \psi_2^*), (s_1^*, s_2^*))$ ,  $\epsilon > 0$ ,  $i \in \{1, 2\}$ , biased belief  $\psi'_i$ , and a post-entry distribution of biased beliefs  $\tilde{\mu}_i \in N_{i,\epsilon}(\psi_i^*, \psi'_i)$ , we say that a post-entry focal configuration  $((\tilde{\mu}_i, \psi_j^*), (s'_i, s'_j))$  is *implausible* if: (1)  $\psi_j^*(s'_i) = \psi_j^*(s_i^*)$ , (2)  $s'_j \neq s_j^*$ , and (3)  $((\tilde{\mu}_i, \psi_j^*), (s'_i, s'_j))$  is a post-entry focal configuration. A post-entry focal configuration is *plausible* if it is not implausible.

**Definition 17.** Configuration  $((\psi_1^*, \psi_2^*), (s_1^*, s_2^*))$  is *plausibly stable* if there exists  $\epsilon > 0$  such that for every player  $i$ , every biased belief  $\psi'_i$ , and every post-entry distribution of biased beliefs  $\tilde{\mu}_i \in N_{i,\epsilon}(\psi_i^*, \psi'_i)$ , there exists a plausible post-entry focal configuration  $((\tilde{\mu}_i, \psi_j^*), (s'_i, s'_j))$  in which the mutants are weakly outperformed relative to the incumbents’ payoff, i.e.,  $\pi_i(s'_i, s'_j) \leq \pi_i(s_i^*, s_j^*)$ .

The following result shows that the definition of a plausibly stable configuration coincides with our definition of BBE. The simple proof, which is analogous to the proof of [10](#), is omitted for brevity.

**Proposition 12.** A configuration  $((\psi_1^*, \psi_2^*), (s_1^*, s_2^*))$  is plausibly stable iff it is a BBE.

### B.3 Illustration of the Evolutionary Interpretation

Similar to DEY, we have presented a reduced-form static notion of evolutionary stability, without formally modeling a detailed dynamics according to which the biased beliefs and the strategies co-evolve. In Section 3.6 we present the essential features of this evolutionary process, which are analogous to DEY’s essential features (see first paragraph in Section 2.2 in DEY): agents are endowed by biased beliefs, these biased-beliefs induce equilibrium behavior in the biased game (presumably by a relatively quick adjustment of the biased players that leads to equilibrium behavior in the biased game), behavior determines “success,” and success (the material payoffs) regulates the evolution of biased beliefs (presumably by a slow process in which agents occasionally die and are replaced by new agents who are more likely to mimic the biased beliefs of more successful incumbents).

In what follows, we illustrate this evolutionary process and its underlying dynamics in an example. Specifically, we present a strong BBE in an “input” game and we illustrate how this strong BBE can persist, given plausible evolutionary dynamics through which the composition of the population evolves.

**Example 13** (*Example 9 revisited*). Consider the following “input” game. Let  $S_i = S_j = [0, 1]$ , and let the payoff function be  $\pi_i(s_i, s_j, \rho) = s_i \cdot s_j - \frac{s_i^2}{2\rho}$ , where the parameter  $\frac{1}{\rho}$  is interpreted as the cost of effort, and we assume that  $\rho \in (0.5, 1)$ . One can show that (1) the best-reply function of each agent is to exert an effort that is  $\rho$  times smaller than the opponent’s (i.e.,  $BR(s_j) = \rho \cdot s_j$ ), (2) in the unique Nash equilibrium of the unbiased game each player exerts no effort  $s_i = s_j = 0$ , and (3) the strategy profile  $(\rho, \rho)$  yields a payoff of  $\rho^2 - \frac{\rho}{2} > 0$ , which is the highest symmetric payoff among all strategy profiles in which agents do not use strictly dominated strategies. Let  $\psi_i^*$  be the following biased-belief function:

$$\psi_i^*(s_j) = \begin{cases} \frac{s_j}{\rho} & s_j < \rho \\ 1 & s_j \geq \rho. \end{cases}$$

In Example 9 we have shown that  $((\psi_1^*, \psi_2^*), (\rho, \rho))$  is a strong BBE. In what follows we illustrate how this strong BBE can persist. Consider a small group of mutants of population  $i$  who have undistorted beliefs. Assume that, initially, the incumbents of population  $j$  use the same strategy against the mutants as they use against the incumbents of population  $i$  (i.e., strategy  $\rho$ ), and the mutants gradually learn to best reply to the incumbents’ behavior by playing  $\rho^2$ . Recall that we assume that the agents of population  $j$  identify the mutants as a separate group of agents who behave differently than the rest of population  $j$  (without assuming that the incumbents of population  $j$  know anything about the biased beliefs of the mutants). These incumbents perceive the mutants’ play as  $\rho$  (due to the incumbents’ biased beliefs), and gradually learn to best reply to this perceived strategy by playing  $\rho^2$ . This, in turn, induces the mutants to adapt their play to playing  $\rho^3$ , and, in response, the incumbents of population  $j$  adapt their play against the mutants and play  $\rho^3$  (the best reply to the mutants’ perceived strategy  $\rho^2$ ). This mutual gradual adaptation process continues until the play in the matches between incumbents of population  $j$  and mutants of population  $i$  converges to  $(0, 0)$ .

Finally, following the convergence of the behavior in the matches against the mutants to  $(0, 0)$ , a slow flow of new agents begins to influence the composition of the population. Each new agent

randomly chooses a mentor among the agents in his own population, where agents with higher fitness are more likely to be chosen as mentors. As the mutants get a much lower payoff (0) than the incumbents of population  $i$  ( $\rho^2 - \frac{\rho}{2} > 0$ ) in the underlying game, their fitness is expected to be lower, and they are much less likely to be chosen as mentors. As a result the share of mutants in the population slowly shrinks until they disappear from the population.

## C Principal-Agent (Subgame-Perfect) Definition of BBE

In this appendix we present an equivalent definition of BBE as a subgame-perfect equilibrium of a two-stage game in which in the first round each player chooses the biased belief of the agent who will play on his behalf in the second round.

### C.1 The Two-Stage Game $\Gamma_G$

Given an underlying two-player normal-form game  $G = (S, \pi)$  define  $\Gamma_G$  as the following four-player two-stage extensive-form game. The four players in the game  $\Gamma$  are: principal 1 and principal 2 (who choose representative agents for the second stage), agent 1 (who plays on behalf of principal 1 in round 2), and agent 2 (who plays on behalf of principal 2 in round 2).

The game  $\Gamma_G$  has 2 stages. In the first stage, the principals simultaneously choose biased beliefs for their agents. That is, each principal  $i$  chooses a biased belief  $\psi_i : S_j \rightarrow S_j$  for agent  $i$ . In the second stage the agents simultaneously choose their strategies. That is, each agent  $i$  chooses strategy  $s_i \in S_i$ . The payoff of each principal  $i$  is  $\pi_i(s_i, s_j)$ . The payoff of each agent  $i$  is  $\pi_i(\psi_i(s_i), s_j)$ . Let  $\Psi_i$  be the set of all feasible (monotone) biased beliefs of agent  $i$ .

A pure strategy profile of  $\Gamma_G$  (henceforth  $\Gamma_G$ -strategy profile) is a tuple  $(\psi_1, \psi_2, \sigma_1, \sigma_2)$ , where each  $\psi_i$  is a biased belief, and each  $\sigma_i : \Psi_1 \times \Psi_2 \rightarrow S_i$  is a function assigning a strategy to each pair of (monotone) biased beliefs. Let  $SPE(\Gamma_G)$  denote the set of all subgame-perfect equilibria of  $\Gamma$ .

### C.2 Subgame-Perfect Definition of Weak BBE

The following result shows that a weak BBE is equivalent to a subgame-perfect equilibrium of  $\Gamma$ . Formally:

**Proposition 13.** *Let  $G$  be a game. Strategy profile  $((\psi_1^*, \psi_2^*), (s_1^*, s_2^*))$  is a weak BBE of  $G$  iff there exists a subgame-perfect equilibrium  $((\psi_1^*, \psi_2^*), (\sigma_1^*, \sigma_2^*))$  of  $\Gamma_G$  satisfying  $\sigma_i^*(\psi_i^*) = s_i^*$  for each player  $i$ .*

*Proof.* “If side”: Let  $((\psi_1^*, \psi_2^*), (\sigma_1^*, \sigma_2^*)) \in SPE(\Gamma_G)$  be a subgame-perfect equilibrium of  $\Gamma$  satisfying  $\sigma_i^*(\psi_i^*) = s_i^*$  for each player  $i$ . Let  $\psi'_i$  be a biased belief of player  $i$ . Let  $s'_1 = \sigma_1^*(\psi'_1, \psi_2^*)$  and  $s'_2 = \sigma_2^*(\psi_1^*, \psi'_2)$ . The fact that  $((\psi_1^*, \psi_2^*), (\sigma_1^*, \sigma_2^*)) \in SPE(\Gamma_G)$  implies that (1)  $(s'_1, s'_2) \in NE\left(G_{(\psi'_1, \psi_2^*)}\right)$  and (2)  $\pi_i(s'_1, s'_2) \leq \pi_i(s_1^*, s_2^*)$ . This implies that  $((\psi_1^*, \psi_2^*), (s_1^*, s_2^*))$  is a weak BBE of  $G$ .

“Only if side”: Let  $((\psi_1^*, \psi_2^*), (s_1^*, s_2^*))$  be a weak BBE of  $G$ . We define  $(\sigma_1^*, \sigma_2^*)$  as follows:<sup>7</sup> (1)  $\sigma_i^*(\psi_1^*, \psi_2^*) = s_i^*$ , (2) for each biased belief  $\psi'_i \neq \psi_i^*$ , define  $\sigma_i^*(\psi'_i, \psi_j^*) = s'_i$  and  $\sigma_j^*(\psi'_i, \psi_j^*) = s'_j$  such that  $(s'_i, s'_j) \in NE(G_{(\psi'_i, \psi_j^*)})$  and  $\pi_i(s'_i, s'_j) \leq \pi_i(s_i^*, s_j^*)$  (such a pair  $(s'_i, s'_j)$  exists due to  $((\psi_1^*, \psi_2^*), (s_1^*, s_2^*))$  being a weak BBE of  $G$ ), and (3) for each pair of biased beliefs  $\psi'_i \neq \psi_i^*$  and  $\psi'_j \neq \psi_j^*$ , define  $\sigma_i^*(\psi'_i, \psi'_j) = s'_i$  and  $\sigma_j^*(\psi'_i, \psi'_j) = s'_j$  such that  $(s'_i, s'_j) \in NE(G_{(\psi'_i, \psi'_j)})$ . The definition of  $(\sigma_1^*, \sigma_2^*)$  immediately implies that  $((\psi_1^*, \psi_2^*), (\sigma_1^*, \sigma_2^*)) \in SPE(\Gamma_G)$ .  $\square$

### C.3 Subgame-Perfect Definition of BBE

Next, we present an equivalent definition of a BBE as a refinement of a subgame-perfect equilibrium of  $\Gamma_G$ . Specifically, a subgame-perfect equilibrium  $(\psi_1^*, \psi_2^*, \sigma_1^*, \sigma_2^*)$  is required to remain a subgame-perfect equilibrium even after changing the off-the-equilibrium path behavior to a different Nash equilibrium of the induced subgame in which (I) a single player (say, player  $j$ ) has deviated to a different biased-belief, (II) the non-deviator perceives the deviator’s strategy in the same way as the original on-the-equilibrium path opponent’s strategy, and (III) the non-deviator changes his behavior such that after the change it coincides with his on-the-equilibrium path behavior. Formally,

**Definition 18.** A subgame-perfect equilibrium  $(\psi_1^*, \psi_2^*, \sigma_1^*, \sigma_2^*) \in SPE(\Gamma_G)$  is a *plausible subgame-perfect equilibrium* if (I) the biased beliefs  $\psi_1^*$  and  $\psi_2^*$  are monotone, and (II)  $(\psi_1^*, \psi_2^*, \sigma'_1, \sigma'_2) \in SPE(\Gamma)$  for each pair of second-stage strategies  $\sigma'_1, \sigma'_2$  satisfying: (1)  $(\psi'_1, \psi'_2, \sigma'_1, \sigma'_2) \in SPE(\Gamma_G)$  for some pair of first-stage strategies  $(\psi'_1, \psi'_2)$  (i.e., second-stage behavior is consistent with equilibrium behavior in all subgames) and (2) if  $\sigma'_i(\psi'_1, \psi'_2) \neq \sigma_i^*(\psi'_1, \psi'_2)$ , then: (I)  $\psi'_i = \psi_i^*$  and  $\psi'_j \neq \psi_j^*$ , (II)  $\psi_i^*(\sigma'_j(\psi'_1, \psi'_2)) = \psi_i^*(\sigma_j^*(\psi'_1, \psi'_2))$ , and (III)  $\sigma'_j(\psi'_1, \psi'_2) = \sigma_j^*(\psi'_1, \psi'_2)$ .

**Proposition 14.** Let  $G$  be a game. Strategy profile  $((\psi_1^*, \psi_2^*), (s_1^*, s_2^*))$  is a BBE of  $G$  iff there exists a plausible subgame-perfect equilibrium  $((\psi_1^*, \psi_2^*), (\sigma_1^*, \sigma_2^*))$  of  $\Gamma_G$  satisfying  $\sigma_i^*(\psi_i^*) = s_i^*$  for each player  $i$ .

The simple proof, which is analogous to the proof of Proposition 13, is omitted for brevity.

### C.4 Subgame-Perfect Definition of Strong BBE

Finally, we present an equivalent definition of a strong BBE as a refinement of a subgame-perfect equilibrium of  $\Gamma$ , which remains an equilibrium even after changing off the equilibrium path in subgames to other Nash equilibria of the induced subgames. Formally,

**Definition 19.** A subgame-perfect equilibrium  $(\psi_1^*, \psi_2^*, \sigma_1^*, \sigma_2^*) \in SPE(\Gamma_G)$  is a *strong subgame-perfect equilibrium* if (I) the biased beliefs  $\psi_1^*$  and  $\psi_2^*$  are monotone, and (II)  $(\psi_1^*, \psi_2^*, \sigma'_1, \sigma'_2) \in SPE(\Gamma_G)$  for each pair of second-stage strategies  $\sigma'_1, \sigma'_2$  satisfying: (1)  $(\psi'_1, \psi'_2, \sigma'_1, \sigma'_2) \in SPE(\Gamma_G)$  for some pair of first-stage strategies  $\psi'_1, \psi'_2$  (i.e., second-stage behavior is consistent with equilibrium behavior in all subgames) and (2)  $\sigma'_i(\psi_1^*, \psi_2^*) = \sigma_i^*(\psi_1^*, \psi_2^*)$  (i.e., behavior after  $(\psi_1^*, \psi_2^*)$  is unchanged).

<sup>7</sup>The definition of  $(\sigma_1^*, \sigma_2^*)$  relies on the axiom of choice.

Our final result shows that a strong BBE is equivalent to a strong subgame-perfect equilibrium of  $\Gamma$ . Formally:

**Proposition 15.** *Let  $G$  be a game. Strategy profile  $((\psi_1^*, \psi_2^*), (s_1^*, s_2^*))$  is a strong BBE of  $G$  iff there exists a strong subgame-perfect equilibrium  $((\psi_1^*, \psi_2^*), (\sigma_1^*, \sigma_2^*))$  of  $\Gamma$  satisfying  $\sigma_i^*(\psi_i^*) = s_i^*$  for each player  $i$ .*

The simple proof, which is analogous to the proof of Proposition 13, is omitted for brevity.

## D Discontinuous Biased Beliefs

In this appendix we present an alternative definition of BBE that relaxes the assumption that biased beliefs have to be continuous. We show that all BBE characterized in the main text remain BBE when deviators are allowed to use discontinuous biased beliefs.

### D.1 Adapted Definitions: Quasi-equilibria

We redefine a *biased belief*  $\psi_i : S_j \rightarrow S_j$  to be an arbitrary (rather than continuous) function that assigns to each strategy of the opponent a (possibly distorted) belief about the opponent's play. The definition of a configuration  $(\psi^*, s^*)$  is left unchanged (i.e., we require that  $(s_i^*, s_j^*) \in NE(G_{\psi^*})$ ).

Recall that a configuration is a BBE if each biased belief is a best reply to the opponent's biased belief, in the sense that an agent who chooses a different biased belief is weakly outperformed in the induced equilibrium of the new biased game. Allowing discontinuous beliefs implies that some biased games  $G_{(\psi_1, \psi_2)}$  in which one (or both) of the biases are discontinuous may not admit Nash equilibria. This requires us to adapt the definition of a BBE to deal with behavior in biased games that do not admit Nash equilibria. We do so by assuming that the resulting behavior in a biased game that does not admit a Nash equilibrium is a “ $j$ -quasi-equilibrium,” in which the non-deviator (player  $j$ ) best replies to the perceived behavior of the deviator (player  $i$ ), while the deviator is allowed to play arbitrarily. Formally:

**Definition 20.** Let  $(\psi_i, \psi_j)$  be a profile of biased beliefs, and let  $j$  be one of the players (interpreted as the non-deviator); then we define  $QE_j(G_{(\psi_i, \psi_j)})$  as the set of  $j$ -quasi-equilibria of the biased game  $G_{(\psi_i, \psi_j)}$  as follows:

$$QE_j(G_{(\psi_i, \psi_j)}) = \begin{cases} NE(G_{(\psi_i, \psi_j)}) & NE(G_{(\psi_i, \psi_j)}) \neq \emptyset \\ \{(s_i, s_j) \mid s_j \in BR(\psi_j(s_i))\} & NE(G_{(\psi_i, \psi_j)}) = \emptyset. \end{cases}$$

Note that any biased game admits a  $j$ -quasi-equilibrium.

### D.2 Adapted Definitions: BBE'

We redefine our notions of BBE as follows, and write them as BBE'. In a strong BBE', the deviator (player  $i$ ) is required to be outperformed in all  $j$ -quasi-equilibria, and biased beliefs are required

to be monotone. In a weak BBE', the deviator is required to be outperformed in at least one  $j$ -quasi-equilibrium. The notion of a BBE' is in between these two notions. Specifically, in a BBE', the biased beliefs are required to be monotone, and, in addition, the deviator (player  $i$ ) is required to be outperformed in at least one plausible  $j$ -quasi-equilibrium of the new biased game, where implausible  $j$ -quasi-equilibria are defined as follows. We say that a  $j$ -quasi-equilibrium of a biased game induced by a deviation of player  $i$  is implausible if (1) player  $i$ 's strategy is perceived by the non-deviating player  $j$  as coinciding with player  $i$ 's original strategy, (2) player  $j$  plays differently relative to his original strategy, and (3) if player  $j$  were playing his original strategy, this would induce a  $j$ -quasi-equilibrium of the biased game. That is, implausible  $j$ -quasi-equilibria are those in which the non-deviating player  $j$  plays differently against a deviator even though player  $j$  has no reason to do so: player  $j$  does not observe any change in player  $i$ 's behavior, and player  $j$ 's original behavior remains an equilibrium of the biased game. Formally:

**Definition 21.** Given configuration  $(\psi^*, s^*)$ , deviating player  $i$ , and biased belief  $\psi'_i$ , we say that a  $j$ -quasi-equilibrium of the biased game  $(s'_i, s'_j) \in QE_j \left( G_{(\psi'_i, \psi_j^*)} \right)$  is *implausible* if: (1)  $\psi_j^*(s'_i) = \psi_j^*(s_i^*)$ , (2)  $s_j^* \neq s'_j$ , and (3)  $(s'_i, s_j^*) \in QE_j \left( G_{(\psi'_i, \psi_j^*)} \right)$ . A  $j$ -quasi-equilibrium is *plausible* if it is not implausible. Let  $PQE_j \left( G_{(\psi'_i, \psi_j^*)} \right)$  be the set of all plausible  $j$ -quasi-equilibria of the biased game  $G_{(\psi'_i, \psi_j^*)}$ .

Note that it is immediate from Definition 21 and the nonemptiness of  $QE_j \left( G_{(\psi'_i, \psi_j^*)} \right)$  that  $PQE_j \left( G_{(\psi'_i, \psi_j^*)} \right)$  is nonempty.

**Definition 22.** Configuration  $(\psi^*, s^*)$  is:

1. a *strong BBE'* if (I) each biased belief  $\psi_i^*$  is monotone, and (II)  $\pi_i(s'_i, s'_j) \leq \pi_i(s_i^*, s_j^*)$  for every player  $i$ , every biased belief  $\psi'_i$ , and every  $j$ -quasi-equilibrium  $(s'_i, s'_j) \in QE_j \left( G_{(\psi'_i, \psi_j^*)} \right)$ ;
2. a *weak BBE'* if for every player  $i$  and every biased belief  $\psi'_i$ , there exists a  $j$ -quasi-equilibrium  $(s'_i, s'_j) \in QE_j \left( G_{(\psi'_i, \psi_j^*)} \right)$ , such that  $\pi_i(s'_i, s'_j) \leq \pi_i(s_i^*, s_j^*)$ ;
3. a *BBE'* if (I) each biased belief  $\psi_i^*$  is monotone, and (II) for every player  $i$  and every biased belief  $\psi'_i$ , there exists a plausible  $j$ -quasi-equilibrium  $(s'_i, s'_j) \in PQE_j \left( G_{(\psi'_i, \psi_j^*)} \right)$ , such that  $\pi_i(s'_i, s'_j) \leq \pi_i(s_i^*, s_j^*)$ .

It is immediate that any strong BBE' is a BBE', and that any BBE' is a weak BBE'.

(resp., strong, weak) BBE'  $(\psi^*, s^*)$  is continuous if each biased function  $\psi_i^*$  is continuous. Note, that deviators are allowed to choose discontinuous biased beliefs.

### D.3 Robustness of BBE to Discontinuous Biased Beliefs

In what follows we observe that all the BBE that we characterize in all the results of the paper are also BBE'. That is, all of our BBE are robust to allowing deviators to use discontinuous biased

beliefs. Specifically, any BBE (resp., weak BBE, strong BBE) that is characterized in any result (or example) in the paper, is a continuous BBE' (resp., weak continuous BBE', strong continuous BBE').

The reason why this observation is true is that in all the arguments in the proofs of the paper's results for why a configuration  $((\psi_i^*, \psi_j^*), (s_i^*, s_j^*))$  is a BBE, when we show that a deviator (player  $i$ ) is outperformed after deviating to biased belief  $\psi_i'$  and after the players play strategy profile  $(s_i', s_j')$ , we rely only on the assumption that the non-deviator (player  $j$ ) best replies to the deviator (i.e., that  $s_j \in BR(\psi_j^*(s_i'))$ ), which is implied by assuming  $(s_i', s_j') \in QE_j(G_{(\psi_i', \psi_j^*)})$ , and we do not use in any of the arguments the assumption that the deviator plays a best reply (i.e., we do not rely on  $s_i \in BR(\psi_i'(s_j'))$  in any of the proofs).

## E Partial Observability

Throughout the paper we assume that if an agent deviates to a different biased belief, then the opponent always observes this deviation. In this appendix, we relax this assumption, and show that our results hold also in a setup with partial observability (some results hold for any level of partial observability, while others hold for a sufficiently high level of observability).

### E.1 Restricted Biased Games

Let  $p \in [0, 1]$  denote the probability that an agent who is matched with an opponent who deviates to a different biased belief *privately* observes the opponent's deviation (henceforth, *observation probability*). If an agent does not observe the deviation, then he continues playing his original configuration's strategy.

Our definitions of configuration and biased game remain unchanged. We now define a restricted biased game  $G_{(\psi_i', \psi_j^*, s_j^*, p)}$  as a game with a payoff function according to which (1) each player's payoff is determined by the opponent's perceived strategy, and (2) the non-deviator is restricted to playing  $s_j^*$  with probability  $p$  (i.e., when not observing the opponent's deviation). Formally:

**Definition 23.** Given an underlying game  $G = (S, \pi)$ , a profile of biased beliefs  $(\psi_i', \psi_j^*)$ , and a strategy  $s_j^*$  of player  $j$  (interpreted as the non-deviator), let the *restricted biased game*  $G_{(\psi_i', \psi_j^*, s_j^*, p)} = (S, \tilde{\pi}(\psi_i', \psi_j^*, s_j^*, p))$  be defined as follows:

$$\tilde{\pi}_i(\psi_i', \psi_j^*, s_j^*, p)(s_i, s_j) = p \cdot \pi_i(s_i, \psi_i'(s_j)) + (1 - p) \cdot \pi_i(s_i, \psi_i'(s_j^*)), \text{ and}$$

$$\tilde{\pi}_j(\psi_i', \psi_j^*, s_j^*, p)(s_i, s_j) = p \cdot \pi_j(s_j, \psi_j^*(s_i)) + (1 - p) \pi_j(s_j^*, \psi_j^*(s_i)).$$

A Nash equilibrium of a  $p$ -restricted biased game is defined in the standard way. Formally, a pair of strategies  $s^* = (s_1', s_2')$  is a Nash equilibrium of a restricted biased game  $G_{(\psi_i', \psi_j^*, s_j^*, p)}$ , if each  $s_i'$  is a best reply against the perceived strategy of the opponent, i.e.,

$$s_i' = \operatorname{argmax}_{s_i \in S_i} (\tilde{\pi}_i(\psi_i', \psi_j^*, s_j^*, p)(s_i, s_j')).$$



Let  $NE\left(G_{(\psi'_i, \psi_j^*, s_j^*)}\right) \subseteq S_1 \times S_2$  denote the set of all Nash equilibria of the restricted biased game  $G_{(\psi'_i, \psi_j^*, s_j^*, p)}$ .

Observe that the set of strategies of a biased game is convex and compact, and the payoff function  $\tilde{\pi}_i(\psi'_i, \psi_j^*, s_j^*, p) : S_i \times S_j \rightarrow \mathbb{R}$  is weakly concave in the first parameter and continuous in both parameters. This implies (due to a standard application of Kakutani's fixed-point theorem) that each restricted biased game  $G_{(\psi'_i, \psi_j^*, s_j^*, p)}$  admits a Nash equilibrium (i.e.,  $NE\left(G_{(\psi'_i, \psi_j^*, s_j^*, p)}\right) \neq \emptyset$ ).

## E.2 $p$ -BBE

We are now ready to define our equilibrium concept. Configuration  $(\psi^*, s^*)$  is a  $p$ -BBE if each biased belief is a best reply to the opponent's biased belief, in the sense that an agent who chooses a different biased belief is weakly outperformed in the induced equilibrium of the new restricted biased game. We present three versions of  $p$ -BBE that differ with respect to the equilibrium selection when the new biased game admits multiple equilibria. In a strong  $p$ -BBE (I) each biased-belief is monotone, and (II) the deviator is required to be outperformed in all Nash equilibria of the new restricted biased game. In a weak BBE, the deviator is required to be outperformed in at least one equilibrium of the new restricted biased game.

The notion of a  $p$ -BBE is in between these two notions. Specifically, in a  $p$ -BBE (I) each biased-belief is monotone, and (II) the deviator is required to be outperformed in at least one plausible equilibrium of the new restricted biased game, where implausible equilibria are defined as follows. We say that a Nash equilibrium of a restricted biased game induced by a deviation of player  $i$  is implausible if (1) player  $i$ 's strategy is perceived by the non-deviating player  $j$  as coinciding with player  $i$ 's original strategy, (2) player  $j$  plays differently relative to his original strategy, and (3) if player  $j$  were playing his original strategy, this would induce an equilibrium of the biased game. That is, implausible equilibria are those in which the non-deviating player  $j$  plays differently against a deviator even though player  $j$  has no reason to do so: player  $j$  does not observe any change in player  $i$ 's behavior, and player  $j$ 's original behavior remains an equilibrium of the biased game. Formally:

**Definition 24.** Given configuration  $(\psi^*, s^*)$ , deviating player  $i$ , and biased belief  $\psi'_i$ , we say that a Nash equilibrium of the restricted biased game  $(s'_i, s'_j) \in NE\left(G_{(\psi'_i, \psi_j^*, s_j^*)}\right)$  is *implausible* if: (1)  $\psi_j^*(s'_i) = \psi_j^*(s^*_i)$ , (2)  $s_j^* \neq s'_j$ , and (3)  $(s'_i, s_j^*) \in NE\left(G_{(\psi'_i, \psi_j^*, s_j^*)}\right)$ . An equilibrium is *plausible* if it is not implausible. Let  $PNE\left(G_{(\psi'_i, \psi_j^*, s_j^*)}\right)$  be the set of all plausible equilibria of the biased game  $G_{(\psi'_i, \psi_j^*, s_j^*, p)}$ .

Note that it is immediate from Definition 24 and the nonemptiness of  $NE\left(G_{(\psi'_i, \psi_j^*, s_j^*)}\right)$  that  $PNE\left(G_{(\psi'_i, \psi_j^*, s_j^*)}\right)$  is nonempty.

**Definition 25.** Configuration  $(\psi^*, s^*)$  is:

1. a *strong p-BBE* if (I) each biased belief  $\psi_i^*$  is monotone, and (II)  $\pi_i(s'_i, s'_j) \leq \pi_i(s_i^*, s_j^*)$  for every player  $i$ , every biased belief  $\psi'_i$ , and every Nash equilibrium  $(s'_i, s'_j) \in NE(G_{(\psi'_i, \psi_j^*, s_j^*, p)})$ ;
2. a *weak p-BBE* if for every player  $i$  and every biased belief  $\psi'_i$ , there exists a Nash equilibrium  $(s'_i, s'_j) \in NE(G_{(\psi'_i, \psi_j^*, s_j^*, p)})$ , such that  $\pi_i(s'_i, s'_j) \leq \pi_i(s_i^*, s_j^*)$ ;
3. a *p-BBE* if (I) each biased belief  $\psi_i^*$  is monotone, and (II) for every player  $i$  and every biased belief  $\psi'_i$ , there exists a plausible Nash equilibrium  $(s'_i, s'_j) \in PNE(G_{(\psi'_i, \psi_j^*, s_j^*, p)})$ , such that  $\pi_i(s'_i, s'_j) \leq \pi_i(s_i^*, s_j^*)$ .

It is immediate that: (1) any strong  $p$ -BBE is a  $p$ -BBE, and that any  $p$ -BBE is a weak  $p$ -BBE, and (2) the definition of 1-BBE (resp., weak 1-BBE, strong 1-BBE) coincides with the original definition of BBE (resp., weak BBE, strong BBE).

### E.3 Extension of Results

In what follows we sketch how to extend our results to the setup of partial observability. The adaptations of the proofs are relatively simple, and, for brevity, we only sketch the differences with respect to the original proofs.

#### E.3.1 Adaptation of Section 4 (Nash Equilibria and BBE Outcomes)

The example that some Nash equilibria cannot be supported as the outcomes of weak  $P$ -BBE with undistorted beliefs can be extended for any  $p > 0$ .

**Example 14** (*Example 1 revisited. Cournot equilibrium cannot be supported by undistorted beliefs*). Consider the following symmetric Cournot game with linear demand  $G = (S, \pi)$ :  $S_i = [0, 1]$  and  $\pi_i(s_i, s_j) = s_i \cdot (1 - s_i - s_j)$  for each player  $i$ . The unique Nash equilibrium of the game is  $s_i^* = s_j^* = \frac{1}{3}$ , which yields both players a payoff of  $\frac{1}{9}$ . Fix observation probability  $p > 0$ . Assume to the contrary that this outcome can be supported as a weak  $p$ -BBE by the undistorted beliefs  $\psi_i^* = \psi_j^* = I_d$ . Fix a sufficiently small  $0 < \epsilon \ll 1$ . Consider a deviation of player 1 to the blind belief  $\psi'_i \equiv \frac{1}{3} - 2 \cdot \epsilon$ . Note that this blind belief has a unique best reply:  $s'_i = \frac{1}{3} + \epsilon$ . The unique equilibrium of the restricted biased game  $G_{(\psi'_i, \psi_j^*, s_j^*, p)}$  is  $s'_j = \frac{1}{3} - \frac{\epsilon}{2}$ ,  $s'_i = \frac{1}{3} + \epsilon$ , which yields the deviator a payoff of  $\frac{1}{9} + \frac{\epsilon}{6} - \frac{\epsilon^2}{2}$  with probability  $p$  (when his deviation is observed by player 2) and a payoff of  $\frac{1}{9} - \epsilon^2$  with probability  $1 - p$  (when his deviation is not observed by player 2). For a sufficiently small  $\epsilon > 0$  the expected payoff of the deviator is strictly larger than  $\frac{1}{9}$ .

All the results of Section 4 hold for any observation probability  $p \in [0, 1]$  with minor adaptations to the proofs.

**Proposition 16** (Proposition 1 extended). *Let  $(s_1^*, s_2^*)$  be a (strict) Nash equilibrium of the game  $G = (S, \pi)$ . Let  $\psi_1^* \equiv s_2^*$  and  $\psi_2^* \equiv s_1^*$ . Then  $((\psi_1^*, \psi_2^*), (s_1^*, s_2^*))$  is a (strong)  $p$ -BBE for any  $p \in [0, 1]$ .*

*Claim 2* (Claim 1 extended). The unique Nash equilibrium payoff of a zero-sum game is also the unique payoff in any weak  $p$ -BBE for any  $p \in [0, 1]$ .

**Proposition 17** (Proposition 2 extended). *If a game admits a strictly dominant strategy  $s_i^*$  for player  $i$ , then any weak  $p$ -BBE outcome is a Nash equilibrium of the underlying game.*

### E.3.2 Adaptation of Section 6 (Main Results)

**Adaptation of Subsection 6.1 (Preliminary Result)** Minor adaptations of the proof of Proposition 3 show that it holds for any  $p \in [0, 1]$ . Formally:

**Proposition 18.** *Let  $p \in [0, 1]$ . If a strategy profile  $s^* = (s_1^*, s_2^*)$  is a weak  $p$ -BBE outcome, then (1) the profile  $s^*$  is undominated and (2)  $\pi_i(s^*) \geq M_i^U$ .*

**Adaptation of Subsection 6.2 (Games with Strategic Complements)** Minor adaptations to the proofs of the results of Subsection 6.2 show that most of these results (namely, part (1) of Proposition 4 and Corollaries 2 and 3) hold for any  $p \in [0, 1]$ , while part (2) of Proposition 4 holds for  $p$ -s sufficiently close to one. Formally:

**Proposition 19** (Proposition 4 extended). *Let  $G$  be a game with strategic substitutes and positive externalities.*

1. *Fix  $p \in [0, 1]$ . Let  $(s_1^*, s_2^*)$  be a  $p$ -BBE outcome. Then  $(s_1^*, s_2^*)$  is (I) undominated, and for each player  $i$ : (II)  $\pi_i(s_i^*, s_j^*) \geq M_i^U$ , and (III)  $s_i^* \leq \max(BR(s_j^*))$  (underinvestment).*
2. *Let  $(s_1^*, s_2^*)$  be an undominated profile satisfying for each player  $i$ : (II')  $\pi_i(s_i^*, s_j^*) > \tilde{M}_i^U$ , and (III)  $s_i^* \leq \max(BR(s_j^*))$ . Then there exists  $\bar{p} < 1$  such that  $(s_1^*, s_2^*)$  is a  $p$ -BBE outcome for any  $p \in [\bar{p}, 1]$ .  
Moreover, if  $\pi_i(s_i, s_j)$  is strictly concave then  $(s_1^*, s_2^*)$  is a strong  $p$ -BBE outcome for any  $p \in [\bar{p}, 1]$ .*

**Corollary 7.** *Fix  $p \in [0, 1]$ . Let  $G$  be a game with strategic complements and positive externalities with a lowest Nash equilibrium  $(\underline{s}_1, \underline{s}_2)$  satisfying  $\underline{s}_1 < \max(S_i)$  for each player  $i$ . Let  $(s_1^*, s_2^*)$  be a  $p$ -BBE outcome. Then  $\underline{s}_i \leq s_i^*$  for each player  $i$ .*

**Corollary 8.** *Fix  $p \in [0, 1]$ . Let  $G$  be a game with positive externalities and strategic complements. Let*

*$((\psi_1^*, \psi_2^*), (s_1^*, s_2^*))$  be a  $p$ -BBE. If  $s_i^* \notin \{\min(S_i), \max(S_i)\}$ , then player  $i$  exhibits wishful thinking (i.e.,  $\psi_i^*(s_j^*) \geq s_j^*$ ).*

One can also adapt the examples of Section 6.2 (and, similarly, the examples of Sections 6.3 and 6.4) to sufficiently high  $p$ s.

**Adaptation of Section 6.3 (Games With Strategic Substitutes)** Minor adaptations to the proofs of the results of Subsection 6.3 show that most of these results (namely, part (1) of Proposition 5 and Corollaries 4 and 5) hold for any  $p \in [0, 1]$ , while part (2) of Proposition 5 holds for  $p$ -s sufficiently close to one. Formally:

**Proposition 20** (Proposition 5 extended). *Let  $G$  be a game with strategic substitutes and positive externalities.*

1. *Fix  $p \in [0, 1]$ . Let  $(s_1^*, s_2^*)$  be a  $p$ -BBE outcome. Then  $(s_1^*, s_2^*)$  is (I) undominated, and for each player  $i$ : (II)  $\pi_i(s_i^*, s_j^*) \geq M_i^U$ , and (III)  $s_i^* \geq \min(BR(s_j^*))$  (overinvestment).*
2. *Let  $(s_1^*, s_2^*)$  be an undominated profile satisfying for each player  $i$ : (II')  $\pi_i(s_i^*, s_j^*) > \tilde{M}_i^U$ , and (III)  $s_i^* \geq \min(BR(s_j^*))$ . Then there exists  $\bar{p} < 1$  such that  $(s_1^*, s_2^*)$  is a  $p$ -BBE outcome for any  $p \in [\bar{p}, 1]$ .  
Moreover, if  $\pi_i(s_i, s_j)$  is strictly concave then  $(s_1^*, s_2^*)$  is a strong  $p$ -BBE outcome for any  $p \in [\bar{p}, 1]$ .*

**Corollary 9.** *Fix  $p \in [0, 1]$ . Let  $G$  be a game with strategic substitutes and positive externalities. Let  $(s_1^*, s_2^*)$  be a BBE outcome. Then, there exists a Nash equilibrium of the underlying game  $(s_1^e, s_2^e)$ , and a player  $i$  such that  $s_i^e \geq s_i^*$ .*

**Corollary 10.** *Fix  $p \in [0, 1]$ . Let  $G$  be a game with strategic substitutes and positive externalities. Let  $((\psi_1^*, \psi_2^*), (s_1^*, s_2^*))$  be a  $p$ -BBE. If  $s_i^* \notin \{\min(S_i), \max(S_i)\}$ , then player  $i$  exhibits wishful thinking (i.e.,  $\psi_i^*(s_j^*) \geq s_j^*$ ).*

**Adaptation of Section 6.3 (Games With Strategic Opposites)** Minor adaptations to the proofs of the results of Subsection 6.3 show that most of these results (namely, part (1) of Proposition 6, and Corollary 6) hold for any  $p \in [0, 1]$ , while part (2) of Proposition 6 holds for  $p$ -s sufficiently close to one. Formally:

**Proposition 21.** *Let  $G$  be a game with positive externalities and strategic opposites:  $\frac{\partial \pi_1(s_1, s_2)}{\partial s_1} > 0$  and  $\frac{\partial \pi_2(s_1, s_2)}{\partial s_1} < 0$  for each pair of strategies  $s_1, s_2$ .*

1. *Fix  $p \in [0, 1]$ . Let  $(s_1^*, s_2^*)$  be a  $p$ -BBE outcome. Then  $(s_1^*, s_2^*)$  is (I) undominated: (II)  $\pi_i(s_i^*, s_j^*) \geq M_i^U$  for each player  $i$ , and (III)  $s_1^* \leq \max(BR(s_2^*))$  and  $s_2^* \geq \min(BR(s_1^*))$  (underinvestment of player 1 and overinvestment of player 2).*
2. *Let  $(s_1^*, s_2^*)$  be a profile satisfying: (I) undominated, (II)  $\pi_i(s_i^*, s_j^*) > \tilde{M}_i^U$  for each player  $i$ , and (III)  $s_1^* \leq \max(BR(s_2^*))$  and  $s_2^* \geq \min(BR(s_1^*))$ . Then there exists  $\bar{p} < 1$  such that  $(s_1^*, s_2^*)$  is a  $p$ -BBE outcome for any  $p \in [\bar{p}, 1]$ .*

**Corollary 11.** *Fix  $p \in [0, 1]$ . Let  $((\psi_1^*, \psi_2^*), (s_1^*, s_2^*))$  be a  $p$ -BBE of a game with positive externalities and strategic opposites (i.e.,  $\frac{\partial \pi_1(s_1, s_2)}{\partial s_1} > 0$  and  $\frac{\partial \pi_2(s_1, s_2)}{\partial s_1} < 0$  for each pair of strategies  $s_1, s_2$ ). If  $s_i^* \notin \{\min(S_i), \max(S_i)\}$ , then player  $i$  exhibits pessimism (i.e.,  $\psi_i^*(s_j^*) \leq s_j^*$ ).*

### E.3.3 Adaptation of Section 7 (Additional Results)

**Adaptation of Subsection 7.1 (BBE with Strategic Stubbornness)** In what follows we show how to extend Example 5 to the setup of partial observability (while we leave the extension of the general result, Proposition 7, to future research). The example focuses on Cournot competition.

We show that for each level of partial observability  $p \in [0, 1]$ , there exists a strong BBE in which one of the players (1) has a blind belief and (2) plays a strategy that is between the Nash equilibrium strategy and the Stackelberg strategy (and the closer it is to the Stackelberg strategy, the higher the value of  $p$ ), while the opponent has undistorted beliefs. The first player's (resp., opponent's) payoff is strictly increasing (resp., decreasing) in  $p$ : it converges to the Nash equilibrium payoff when  $p \rightarrow 0$ , and it converges to the Stackelberg leader's (resp., follower's) payoff when  $p \rightarrow 1$ .

**Example 15** (*Example 5 revisited*). Consider the symmetric Cournot game with linear demand:  $G = (S, \pi)$ :  $S_i = \mathbb{R}^+$  and  $\pi_i(s_i, s_j) = s_i \cdot (1 - s_i - s_j)$  for each player  $i$ . Let  $p \in [0, 1]$  be the observation probability. Then

$$\left( \left( \frac{1-p}{3-p}, I_d \right), \left( \frac{1}{3-p}, \frac{2-p}{2 \cdot (3-p)} \right) \right)$$

is a strong BBE that yields a payoff of  $\frac{2-p}{2 \cdot (3-p)}$  to player 1, and yields a payoff of  $\left( \frac{2-p}{2 \cdot (3-p)} \right)^2$  to player 2. Observe that player 1's payoff is increasing in  $p$ , and it converges to the Nash equilibrium (resp., Stackelberg leader's) payoff of  $\frac{1}{9}$  ( $\frac{1}{8}$ ) when  $p \rightarrow 0$  ( $p \rightarrow 1$ ). Further observe that player 2's payoff is decreasing in  $p$ , and it converges to the Nash equilibrium (resp., Stackelberg follower's) payoff of  $\frac{1}{9}$  ( $\frac{1}{16}$ ) when  $p \rightarrow 0$  ( $p \rightarrow 1$ ). The argument that  $\left( \left( \frac{1-p}{3-p}, I_d \right), \left( \frac{1}{3-p}, \frac{2-p}{2 \cdot (3-p)} \right) \right)$  is a strong BBE is sketched as follows: (1)  $\left\{ \left( \frac{1-p}{3-p}, \frac{2-p}{2 \cdot (3-p)} \right) \right\} = NE \left( G_{\left( \frac{1-p}{3-p}, I_d \right)} \right)$  (because  $\frac{1}{3-p}$  is the unique best reply against  $\frac{1-p}{3-p}$  and  $\frac{2-p}{2 \cdot (3-p)}$  is the unique best reply against  $\frac{1}{3-p}$ ); (2) for any biased belief  $\psi'_2$ , player 1 keeps playing  $\frac{1}{3-p}$  due to having a blind belief, and as a result player 2's payoff is at most  $\left( \frac{2-p}{2 \cdot (3-p)} \right)^2$ ; and (3) for any biased belief  $\psi'_1$  inducing a deviating player 1 to play strategy  $x$ , player 2 plays  $\frac{1-x}{2}$  (the unique best reply against  $x$ ) with probability  $p$  (when observing the deviation), and player 2 plays  $\frac{2-p}{2 \cdot (3-p)}$  (the original configuration strategy) with the remaining probability of  $1-p$ . Thus, the payoff of a deviating player 1 who deviates into playing strategy  $x$  is

$$\pi(x) := p \cdot x \cdot \left( \frac{1-x}{2} \right) + (1-p) \cdot x \cdot \left( 1-x - \frac{2-p}{2 \cdot (3-p)} \right) = \left( 1 - \frac{p}{2} \right) \cdot x \cdot (1-x) - \frac{(2-p) \cdot (1-p)}{2 \cdot (3-p)} \cdot x,$$

where this payoff function  $\pi(x)$  is strictly concave in  $x$  with a unique maximum at  $x = \frac{1}{3-p}$  (the unique solution to the FOC  $0 = \frac{\partial \pi}{\partial x} = (1 - \frac{p}{2}) \cdot (1 - 2 \cdot x) - \frac{(2-p) \cdot (1-p)}{2 \cdot (3-p)}$ ).

**Extending the Folk Theorem Results for Sufficiently High  $p$ -s** The main results of Subsection 7.2, show folk theorem results for: (1) monotone BBE in games that admit best replies with full undominated support, and (2) strong BBE in interval games with a payoff function that is strictly concave in the agent's strategy, and weakly convex in the opponent's strategy. Minor adaptations of each proof can show that each result can be extended to  $p$ -s that are sufficiently close to one. Formally:

**Proposition 22** (*Proposition 8 extended*). Let  $G$  be a finite game that admits best replies with full undominated support. Let  $(s_1^*, s_2^*)$  be an undominated strategy profile that induces for each player a payoff above his minmax payoff (i.e.,  $\pi_i(s_1^*, s_2^*) > M_i^U \ \forall i \in \{1, 2\}$ ). Then there exists  $\bar{p} < 1$ , such that  $(s_1^*, s_2^*)$  is a monotone weak  $p$ -BBE outcome for each  $p \in [\bar{p}, 1]$ .

**Proposition 23** (*Proposition 9 extended*). *Let  $G = (S, \pi)$  be an interval game. Assume that for each player  $i$ ,  $\pi_i(s_i, s_j)$  is strictly concave in  $s_i$  and weakly convex in  $s_j$ . If  $(s_1^*, s_2^*)$  is undominated and  $\pi_i(s_1^*, s_2^*) > M_i^U$  for each player  $i$ , then there exists  $\bar{p} < 1$ , such that  $(s_1^*, s_2^*)$  is a strong  $p$ -BBE outcome for each  $p \in [\bar{p}, 1]$ .*

*Sketch of adapting the proofs of Propositions 22 and 23 to the setup of partial observability.* Observe that the gain of an agent who deviates to a different biased belief, when his deviation is unobserved by the opponent, is bounded (due to the payoff of the underlying game being bounded). When the deviation is observed by the opponent, the agent is strictly outperformed, given the BBE constructed in the proofs of Propositions 8 and 9. This implies that there exists  $\bar{p} < 1$  sufficiently close to one, such that the loss of a mutant when being observed by his opponent outweighs the mutant's gain when being unobserved for any  $p \in [\bar{p}, 1]$ .  $\square$

## F Proofs

### F.1 Proof of Proposition 9

Recall that we assume the payoff function  $\pi_i$  to be continuously twice differentiable. This implies that  $\pi_i$  is Lipschitz continuous. Let  $K_i > 0$  be the Lipschitz constant of the payoff function  $\pi_i$  with respect to its first parameter, i.e.,  $K_i$  satisfies

$$\|\pi_i(s_1, s_2) - \pi_i(s'_1, s_2)\| \leq K_i \cdot \|s_1 - s'_1\|.$$

Assume that  $(s_1^*, s_2^*)$  is undominated and  $\pi_i(s_1^*, s_2^*) > M_i^U$  for each player  $i$ . Let  $0 < D_i = \pi_i(s_1^*, s_2^*) - M_i^U$ . For each player  $j$ , let  $s_j^p$  be an undominated strategy that guarantees that player  $i$  obtains, at most, his minmax payoff  $M_i^U$ , i.e.,  $s_j^p = \operatorname{argmin}_{s_j \in S_j^U} (\max_{s_i \in S_i} \pi_i(s_i, s_j))$ . The strict concavity of  $\pi_i(s_i, s_j)$  with respect to  $s_i$  implies that the best-reply correspondence is a continuous one-to-one function. Thus,  $BR^{-1}(s_i)$  is a singleton for each  $s_i$ , and we identify  $BR^{-1}(s_i)$  with the unique element in this singleton set.

Let  $\epsilon > 0$  be a sufficiently small number satisfying  $\epsilon < \min\left(\frac{D_i}{K_i}, \frac{D_j}{K_j}\right)$ . For each  $\delta \in [0, 1]$  define for each player  $i$ :

$$s_i^\delta = \frac{\epsilon - \delta}{\epsilon} \cdot s_i^* + \frac{\delta}{\epsilon} \cdot s_i^p.$$

Let  $\psi_i^\epsilon$  be defined as follows:

$$\psi_i^\epsilon(s'_j) = \begin{cases} BR^{-1}\left(s_i^{|s'_j - s_j|}\right) & |s'_j - s_j| < \epsilon \\ BR^{-1}(s_i^p) & |s'_j - s_j| \geq \epsilon. \end{cases}$$

Note that  $\psi_i^\epsilon$  is continuous. We now show that  $((\psi_1^\epsilon, \psi_2^\epsilon), (s_1^*, s_2^*))$  is a strong BBE. Observe first that the definition of  $(\psi_1^\epsilon, \psi_2^\epsilon)$  immediately implies that  $(s_1^*, s_2^*) \in NE\left(G_{(\psi_1^\epsilon, \psi_2^\epsilon)}\right)$ . Next, consider a deviation of player  $i$  to an arbitrary biased belief  $\psi'_i$ . Consider any equilibrium of the new biased game  $(s'_i, s'_j) \in NE\left(G_{(\psi'_i, \psi_j^\epsilon)}\right)$ . If  $|s'_i - s_i| \geq \epsilon$ , then the definition of  $\psi_j^\epsilon(s'_i)$  implies that  $s_j^p = s'_j$ , and that player  $i$  achieves a payoff of at most  $M_i^U < \pi_i(s_1^*, s_2^*)$ . If  $s'_i = s_i^*$ , then it is immediate that

$s'_j = s_j^*$  and that player  $i$  does not gain from his deviation. If  $0 < |s'_i - s_i| < \epsilon$ , then the definition of  $\psi_j^\epsilon(s'_i)$  implies that

$$\begin{aligned}
\pi_i(s'_i, s'_j) &= \pi_i\left(s'_i, s_j \left| \frac{s'_i - s_i}{\epsilon} \right| \right) = \pi_i\left(s'_i, \frac{\epsilon - |s'_i - s_i|}{\epsilon} \cdot s_j^* + \frac{|s'_i - s_i|}{\epsilon} \cdot s_j^p\right) \leq \\
&\frac{\epsilon - |s'_i - s_i|}{\epsilon} \cdot \pi_i(s'_i, s_j^*) + \frac{|s'_i - s_i|}{\epsilon} \cdot \pi_i(s'_i, s_j^p) \leq \frac{\epsilon - |s'_i - s_i|}{\epsilon} \cdot \pi_i(s'_i, s_j^*) + \frac{|s'_i - s_i|}{\epsilon} \cdot M_i^U \leq \\
&\frac{\epsilon - |s'_i - s_i|}{\epsilon} \cdot \pi_i(s_i^*, s_j^*) + K_i \cdot |s'_i - s_i| + \frac{|s'_i - s_i|}{\epsilon} \cdot M_i^U = \\
&\frac{\epsilon - |s'_i - s_i|}{\epsilon} \cdot \pi_i(s_i^*, s_j^*) + K_i \cdot |s'_i - s_i| + \frac{|s'_i - s_i|}{\epsilon} \cdot (\pi_i(s_i^*, s_j^*) - D_i) = \\
\pi_i(s_i^*, s_j^*) &+ \frac{\epsilon - |s'_i - s_i|}{\epsilon} \cdot K_i \cdot |s'_i - s_i| - \frac{|s'_i - s_i|}{\epsilon} \cdot D_i \leq \pi_i(s_i^*, s_j^*) + K_i \cdot |s'_i - s_i| - \frac{|s'_i - s_i|}{\epsilon} \cdot D_i = \\
&\pi_i(s_i^*, s_j^*) + |s'_i - s_i| \cdot \left(K_i - \frac{D_i}{\epsilon}\right) < \pi_i(s_i^*, s_j^*),
\end{aligned}$$

where the first inequality is due to the convexity of  $\pi_i(s_i, s_j)$  with respect to  $s_j$ , the second inequality is due to  $\pi_i(s'_i, s_j^p) \leq M_i^U$ , the third inequality is due to the Lipschitz continuity, the penultimate inequality is implied by  $\frac{\epsilon - |s'_i - s_i|}{\epsilon} < 1$ , and the last inequality is due to defining  $\epsilon$  to satisfy  $\epsilon < \min\left(\frac{D_i}{K_i}, \frac{D_j}{K_j}\right)$ . This proves that player  $i$  cannot gain from his deviation, and that  $((\psi_1^\epsilon, \psi_2^\epsilon), (s_1, s_2))$  is a strong BBE.

## F.2 Proof of Proposition 4

**Part 1:** Proposition 3 implies (I) and (II). It remains to show (III, overinvestment). Let  $((\psi_i^*, \psi_j^*), (s_i^*, s_j^*))$  be a BBE. Assume to the contrary that  $s_i^* < \min(BR(s_j^*))$ . Consider a deviation of player  $i$  to a blind belief that the opponent always plays strategy  $s_j^*$  (i.e.,  $\psi_i' \equiv s_j^*$ ). Let  $(s'_i, s'_j) \in PNE\left(G_{(\psi_i', \psi_j^*)}\right)$  be a plausible equilibrium of the new biased game. Observe first that  $s'_i \in BR(\psi_i'(s'_j)) = BR(s_j^*)$ . This implies that  $s'_i > s_i^*$ , and, thus, due to the monotonicity of  $\psi_j^*$  we have:  $\psi_j^*(s'_i) \geq \psi_j^*(s_i^*)$ . We consider two cases:

1. If  $\psi_j^*(s'_i) > \psi_j^*(s_i^*)$ , then the strategic complementarity implies that  $s'_j \geq \min(BR(\psi_j^*(s'_i))) \geq \max(BR(\psi_j^*(s_i^*))) \geq s_j^*$ , and this, in turn, implies that player  $i$  strictly gains from his deviation:  $\pi_i(s'_i, s'_j) \geq \pi_i(s'_i, s_j^*) > \pi_i(s_i^*, s_j^*)$ , a contradiction.
2. If  $\psi_j^*(s'_i) = \psi_j^*(s_i^*)$ , then  $(s'_i, s_j^*) \in PNE\left(G_{(\psi_i', \psi_j^*)}\right)$  and  $\pi_i(s'_i, s_j^*) > \pi_i(s_i^*, s_j^*)$ , which contradicts that  $((\psi_i^*, \psi_j^*), (s_i^*, s_j^*))$  is a BBE.

**Part 2:** Assume that strategy profile  $(s_1^*, s_2^*)$  satisfies I, II, and III. For each player  $i$  let  $s_i^\epsilon = \min(BR^{-1}(s_i^*))$ . For every player  $i$  and every strategy  $s_i < s_i^*$  define  $X(s_i)$  as the set of strategies  $s'_i$  for which player  $i$  is worse off (relative to  $\pi_i(s_1^*, s_2^*)$ ) if he plays strategy  $s_i$ , while player  $j$  plays



a best-reply to  $s'_i$ . Formally:

$$X_{s^*}(s_i) = \left\{ s'_i \in S_i \mid \pi_i(s_i, s_j) \leq \pi_i(s_i^*, s_j^*) \quad \forall s_j \in BR(s'_i) \right\}.$$

The assumption that  $\pi_i(s_i^*, s_j^*) > \tilde{M}_i^U$  implies that  $X_{s^*}(s_i)$  is nonempty for each  $s_i$ . The assumption of strategic complements implies that  $X_{s^*}(s_i)$  is an interval starting at  $\min(S_i)$ . Let  $\phi_{s^*}(s_i) = \sup(X_{s^*}(s_i))$ . The assumption that the payoff function is continuously twice differentiable implies that  $\phi_{s^*}(s_i)$  is continuous. The assumption that  $s_j^e = \min(BR^{-1}(s_i^*))$  implies that  $\lim_{s_i \nearrow s_i^*} (\phi_{s^*}(s_i)) = s_i^e$ . These observations imply that for each player  $j$  there exists a monotone biased belief  $\psi_j^*$  satisfying (1)  $\psi_j^*(s_i) = s_i^e$  for each  $s_i \geq s_i^*$  and (2)  $\psi_j^*(s_i) \leq \phi_{s^*}(s_i)$  for each  $s_i < s_i^*$  with an equality only if  $\phi_{s^*}(s_i) = \min(S_i)$ .

We now show that these properties of  $(\psi_1^*, \psi_2^*)$  imply that  $((\psi_1^*, \psi_2^*), (s_1^*, s_2^*))$  is a BBE (a strong BBE if  $\pi_i(s_i, s_j)$  is strictly concave in  $s_i$ ). Consider a deviation of player  $i$  into an arbitrary biased belief  $\psi'_i$ . For each  $s'_i \geq s_i^*$ , and each  $(s'_i, s'_j) \in PNE\left(G_{(\psi'_i, \psi_j^*)}\right)$  ( $(s'_i, s'_j) \in NE\left(G_{(\psi'_i, \psi_j^*)}\right)$ ), the fact that  $\psi_j^*(s'_i) = \psi_j^*(s_i^*)$  implies that  $s'_j = s_j^*$ , and due to assumption (III) of overinvestment and the concavity of the payoff function:  $\pi_i(s'_i, s'_j) = \pi_i(s'_i, s_j^*) \leq \pi_i(s_i^*, s_j^*)$ . For each  $s'_i < s_i^*$ , and each  $(s'_i, s'_j) \in NE\left(G_{(\psi'_i, \psi_j^*)}\right)$ , the fact that  $\psi_j^*(s'_i) \leq \phi_{s^*}(s'_i)$  with an equality only if  $\phi_{s^*}(s'_i) = \min(S_i)$  (and, thus,  $\psi_j^*(s'_i) \in X_{s^*}(s'_i)$ ) implies that  $\pi_i(s'_i, s'_j) \leq \pi_i(s_i^*, s_2^*)$ . This shows that player  $i$  cannot gain from his deviation, which implies that  $((\psi_1^*, \psi_2^*), (s_1^*, s_2^*))$  is a (strong) BBE.

### F.3 Proof of a Lemma Required for Corollary 2

**Lemma 1.** *Let  $G$  be a game with strategic complements and positive externalities with a lowest Nash equilibrium  $(\underline{s}_1, \underline{s}_2)$  satisfying  $\underline{s}_i < \max(S_i)_i$  for each player  $i$ . Let  $s_1^* < \underline{s}_1$ . Then for each  $s_2^* \in S_2$  either (1)  $s_1^* < \min(BR(s_2^*))$  or (2)  $s_2^* < \min(BR(s_1^*))$ .*

*Proof.* Assume first that  $s_2^* > \underline{s}_2$ . The fact that  $\underline{s}_1 \in BR(\underline{s}_2)$  and  $s_2^* > \underline{s}_2$ , together with the strategic complements, imply that  $s_1^* < \underline{s}_1 < \min(BR(s_2^*))$ . We are left with the case where  $s_2^* \leq \underline{s}_2$ . Consider a restricted game in which the set of strategies of each player  $i$  is restricted to being strategies that are at most  $s_i^*$ . The game is a game of strategic complements, and, thus, it admits a pure Nash equilibrium  $(s'_i, s'_j)$ . The minimality of  $(\underline{s}_1, \underline{s}_2)$  implies that  $(s'_i, s'_j)$  cannot be a Nash equilibrium of the unrestricted game. The strategic complements and the concavity of the payoff jointly imply that if  $(s'_i, s'_j)$  is not a Nash equilibrium of the unrestricted game, then there is player  $i$  for which  $s_i^* = s'_i < \min(BR(s'_j)) \leq \min(BR(s_j^*))$ .  $\square$

### F.4 Proof of a Lemma Required for Corollary 3

**Lemma 2.** *Let  $G$  be a game with positive externalities and strategic complementarity of the payoff of player  $i$  (i.e.,  $\frac{\partial^2 \pi_i(s_i, s_j)}{\partial s_i \partial s_j} > 0$  for each  $s_i, s_j$ ). Then  $s'_j < s_j$  implies that  $\max(BR(s'_j)) \leq \min(BR(s_j))$  with an equality only if  $\max(BR(s'_j)) = \min(BR(s_j)) \in \{\min(S_i), \max(S_i)\}$ .*

*Proof.* The inequality  $s'_j < s_j$  and the strategic complementarity of the payoff of player  $i$  implies that  $\frac{\partial \pi_i(s_i, s'_j)}{\partial s_i} < \frac{\partial \pi_i(s_i, s_j)}{\partial s_i}$  for each  $s_i \in S_i$ , which implies that whenever  $\max(BR(s'_j)) \notin$

$\{\min(S_i), \max(S_i)\}$ , then

$$\begin{aligned} \max(BR(s'_j)) &= \max \left\{ s_i^* \in S_i \mid \frac{\partial \pi_i(s_i, s'_j)}{\partial s_i} = 0 \Big|_{s_i=s_i^*} \right\} \\ &< \min \left\{ s_i^* \in S_i \mid \frac{\partial \pi_i(s_i, s_j)}{\partial s_i} = 0 \Big|_{s_i=s_i^*} \right\} \leq \min(BR(s_j)). \end{aligned}$$

This shows that the strict inequality holds whenever  $\max(BR(s'_j)) \notin \{\min(S_i), \max(S_i)\}$ . It remains to show that the weak inequality (namely,  $\max(BR(s'_j)) \leq \min(BR(s_j))$ ) holds when  $\max(BR(s'_j)) \in \{\min(S_i), \max(S_i)\}$ . If  $\max(BR(s'_j)) = \min(S_i)$  then this is immediate. Assume that  $\max(BR(s'_j)) = \max(S_i)$ . Then:

$$\begin{aligned} \max(S_i) &= \max \left\{ s_i^* \in S_i \mid \frac{\partial \pi_i(s_i, s'_j)}{\partial s_i} \geq 0 \Big|_{s_i=s_i^*} \right\} \\ &\leq \min \left\{ s_i^* \in S_i \mid \frac{\partial \pi_i(s_i, s_j)}{\partial s_i} \geq 0 \Big|_{s_i=s_i^*} \right\} \leq \min(BR(s_j)). \end{aligned}$$

□

## F.5 Proof of Proposition 5

The proof is analogous to the proof of Proposition 4, and is presented for completeness.

**Part 1:** Proposition 3 implies (I) and (II). It remains to show (III) (underinvestment). Let  $((\psi_i^*, \psi_j^*), (s_i^*, s_j^*))$  be a BBE. Assume to the contrary that  $s_i^* > \max(BR(s_j^*))$ . Consider a deviation of player  $i$  to a blind belief that the opponent always plays strategy  $s_j^*$  (i.e.,  $\psi'_i \equiv s_j^*$ ). Let  $(s'_i, s'_j) \in PNE(G_{(\psi'_i, \psi_j^*)})$  be a plausible equilibrium of the new biased game. Observe first that  $s'_i \in BR(\psi'_i(s'_j)) = BR(s_j^*)$ . This implies that  $s'_i < s_i^*$ , and, thus, due to the monotonicity of  $\psi_j^*$  we have:  $\psi_j^*(s'_i) \leq \psi_j^*(s_i^*)$ . We consider two cases:

1. If  $\psi_j^*(s'_i) < \psi_j^*(s_i^*)$ , then the strategic substitutability implies that  $s'_j \geq \min(BR(\psi_j^*(s'_i))) \geq \max(BR(\psi_j^*(s_i^*))) \geq s_j^*$ , and this, in turn, implies that player  $i$  strictly gains from his deviation:  $\pi_i(s'_i, s'_j) \geq \pi_i(s'_i, s_j^*) > \pi_i(s_i^*, s_j^*)$ , a contradiction.
2. If  $\psi_j^*(s'_i) = \psi_j^*(s_i^*)$ , then  $(s'_i, s_j^*) \in PNE(G_{(\psi'_i, \psi_j^*)})$  and  $\pi_i(s'_i, s_j^*) > \pi_i(s_i^*, s_j^*)$ , which contradicts that  $((\psi_i^*, \psi_j^*), (s_i^*, s_j^*))$  is a BBE.

**Part 2:** Assume that strategy profile  $(s_1^*, s_2^*)$  satisfies I, II, and III. For each player  $i$  let  $s_i^e = \max(BR^{-1}(s_i^*))$ . For each player  $i$  and each strategy  $s_i > s_i^e$  define  $X(s_i)$  as the set of strategies  $s'_i$  for which player  $i$  is worse off (relative to  $\pi_i(s_1^*, s_2^*)$ ) if he plays strategy  $s_i$ , while player  $j$  plays

a best-reply to  $s'_i$ . Formally:

$$X_{s^*}(s_i) = \left\{ s'_i \in S_i \mid \pi_i(s_i, s_j) \leq \pi_i(s_i^*, s_j^*) \quad \forall s_j \in BR(s'_i) \right\}.$$

The assumption that  $\pi_i(s_i^*, s_j^*) > \tilde{M}_i^U$  implies that  $X_{s^*}(s_i)$  is nonempty for each  $s_i$ . The assumption of strategic substitutes implies that  $X_{s^*}(s_i)$  is an interval ending at  $\max(S_i)$ . Let  $\phi_{s^*}(s_i) = \inf(X_{s^*}(s_i))$ . The assumption that the payoff function is continuously twice differentiable implies that  $\phi_{s^*}(s_i)$  is continuous. The assumption that  $s_j^e = \max(BR^{-1}(s_i^*))$  implies that  $\lim_{s_i \searrow s_i^*} (\phi_{s^*}(s_i)) = s_i^e$ . These observations imply that for each player  $j$  there exists a monotone biased belief  $\psi_j^*$  satisfying (1)  $\psi_j^*(s_i) = s_i^e$  for each  $s_i \leq s_i^*$  and (2)  $\psi_j^*(s_i) \geq \phi_{s^*}(s_i)$  for each  $s_i > s_i^*$  with an equality only if  $\phi_{s^*}(s_i) = \max(S_i)$ .

We now show that these properties of  $(\psi_1^*, \psi_2^*)$  imply that  $((\psi_1^*, \psi_2^*), (s_1^*, s_2^*))$  is a BBE (a strong BBE if  $\pi_i(s_i, s_j)$  is strictly concave in  $s_i$ ). Consider a deviation of player  $i$  to an arbitrary biased belief  $\psi'_i$ . For each  $s'_i \leq s_i^*$ , and each  $(s'_i, s'_j) \in PNE\left(G_{(\psi'_i, \psi_j^*)}\right)$  ( $(s'_i, s'_j) \in NE\left(G_{(\psi'_i, \psi_j^*)}\right)$ ), the fact that  $\psi_j^*(s'_i) = \psi_j^*(s_i^*)$  implies that  $s'_j = s_j^*$  and, due to assumption (III) of underinvestment and the concavity of the payoff function, it follows that  $\pi_i(s'_i, s'_j) = \pi_i(s'_i, s_j^*) \leq \pi_i(s_i^*, s_j^*)$ . For each  $s'_i > s_i^*$ , and each  $(s'_i, s'_j) \in NE\left(G_{(\psi'_i, \psi_j^*)}\right)$ , the fact that  $\psi_j^*(s'_i) \geq \phi_{s^*}(s'_i)$  with an equality only if  $\phi_{s^*}(s_i) = \max(S_i)$  (and, thus,  $\psi_j^*(s'_i) \in X_{s^*}(s'_i)$ ) implies that  $\pi_i(s'_i, s'_j) \leq \pi_i(s_1^*, s_2^*)$ . This shows that player  $i$  cannot gain from his deviation, which implies that  $((\psi_1^*, \psi_2^*), (s_1^*, s_2^*))$  is a (strong) BBE.

## F.6 Proof of a Lemma Required for Corollary 4

**Lemma 3.** *Let  $G$  be a game with strategic substitutes and positive externalities. Let  $(s_1^*, s_2^*)$  be a strategy profile satisfying  $s_i^* > s_i^e$  for each player  $i$  and each Nash equilibrium  $(s_1^e, s_2^e) \in NE(G)$ . Then, either (1)  $s_1^* > \max(BR(s_2^*))$  or (2)  $s_2^* > \max(BR(s_1^*))$ .*

*Proof.* Consider a restricted game in which the set of strategies of each player  $i$  is restricted to being strategies that are at least  $s_i^*$ . The restricted game is a game with strategic substitutes, and, thus, it admits a pure Nash equilibrium  $(s'_1, s'_2)$  (recall, that after relabeling the set of strategies of one of the players, the game becomes supermodular, and because of this the game admits a pure Nash equilibrium due to [Milgrom and Roberts, 1990](#)). The assumption that  $s_i^* > s_i^e$  for each player  $i$  and each Nash equilibrium  $(s_1^e, s_2^e) \in NE(G)$  implies that  $(s'_1, s'_2)$  cannot be a Nash equilibrium of the unrestricted game. The concavity of the payoff and the strategic substitutes jointly imply that if  $(s'_i, s'_j)$  is not a Nash equilibrium of the unrestricted game, then there is a player  $i$  for which  $s_i^* = s'_i > \max(BR(s'_j)) \geq \max(BR(s_j^*))$ .  $\square$

## F.7 Proof of Corollary 4

The proof is analogous to Corollary 3, and is presented for completeness. Assume to the contrary that  $\psi_i^*(s_j^*) < s_j^*$ . Lemma 4 (below) implies that  $\min(BR(\psi_i^*(s_j^*))) \geq \max(BR(s_j^*))$  with an

equality only if

$$\min \left( BR \left( \psi_i^* \left( s_j^* \right) \right) \right) \in \{ \min (S_i), \max (S_i) \}.$$

Part 1 of Proposition 5 and the definition of a monotone BBE imply that

$$\min \left( BR \left( \psi_i^* \left( s_j^* \right) \right) \right) \leq s_i^* \leq \max \left( BR \left( s_j^* \right) \right).$$

The previous inequalities jointly imply that

$$\min \left( BR \left( \psi_i^* \left( s_j^* \right) \right) \right) = s_i^* = \max \left( BR \left( s_j^* \right) \right) \in \{ \min (S_i), \max (S_i) \},$$

which contradicts the assumption that  $s_i^* \notin \{ \min (S_i), \max (S_i) \}$ .

**Lemma 4.** *Let  $G$  be a game with positive externalities and strategic substitutability of the payoff of player  $i$  (i.e.,  $\frac{\partial^2 \pi_i(s_i, s_j)}{\partial s_i \partial s_j} > 0$  for each  $s_i, s_j$ ). Then  $s'_j < s_j$  implies that  $\min \left( BR \left( s'_j \right) \right) \geq \max \left( BR \left( s_j \right) \right)$  with an equality only if  $\min \left( BR \left( s'_j \right) \right) = \min \left( BR \left( s_j \right) \right) \in \{ \min (S_i), \max (S_i) \}$ .*

*Proof.* The proof is analogous to the proof of Lemma 1, and is presented for completeness. The inequality  $s'_j < s_j$  and the strategic substitutability of the payoff of player  $i$  implies that  $\frac{\partial \pi_i(s_i, s'_j)}{\partial s_i} > \frac{\partial \pi_i(s_i, s_j)}{\partial s_i}$  for each  $s_i \in S_i$ , which implies that whenever  $\min \left( BR \left( s'_j \right) \right) \notin \{ \min (S_i), \max (S_i) \}$ , then

$$\begin{aligned} \min \left( BR \left( s'_j \right) \right) &= \min \left\{ s_i^* \in S_i \mid \frac{\partial \pi_i(s_i, s'_j)}{\partial s_i} = 0 \mid_{s_i = s_i^*} \right\} \\ &> \max \left\{ s_i^* \in S_i \mid \frac{\partial \pi_i(s_i, s_j)}{\partial s_i} = 0 \mid_{s_i = s_i^*} \right\} = \max \left( BR \left( s_j \right) \right). \end{aligned}$$

This shows that the strict inequality holds whenever  $\min \left( BR \left( s'_j \right) \right) \notin \{ \min (S_i), \max (S_i) \}$ . It remains to show that the weak inequality (namely,  $\min \left( BR \left( s'_j \right) \right) \geq \max \left( BR \left( s_j \right) \right)$ ) holds when  $\min \left( BR \left( s'_j \right) \right) \in \{ \min (S_i), \max (S_i) \}$ . If  $\min \left( BR \left( s'_j \right) \right) = \max (S_i)$  then this is immediate. Assume that  $\min \left( BR \left( s'_j \right) \right) = \min (S_i)$ . Then:

$$\begin{aligned} \min (S_i) &= \min \left\{ s_i^* \in S_i \mid \frac{\partial \pi_i(s_i, s'_j)}{\partial s_i} \geq 0 \mid_{s_i = s_i^*} \right\} \\ &\geq \max \left\{ s_i^* \in S_i \mid \frac{\partial \pi_i(s_i, s_j)}{\partial s_i} \geq 0 \mid_{s_i = s_i^*} \right\} \geq \max \left( BR \left( s_j \right) \right). \end{aligned}$$

□

## F.8 Proof of Proposition 6

The proof is analogous to the proof of Proposition 4, and is presented for completeness.

**Part 1:** Proposition 3 implies (I) and (II). It remains to show (III). Let  $((\psi_1^*, \psi_2^*), (s_1^*, s_2^*))$  be a BBE. We begin by showing overinvestment of player 2. Assume to the contrary that  $s_2^* <$

$\min(BR(s_j^*))$ . Consider a deviation of player 2 to a blind belief that the opponent always plays strategy  $s_1^*$  (i.e.,  $\psi_2' \equiv s_1^*$ ). Let  $(s_1', s_2') \in PNE(G_{(\psi_1^*, \psi_2')})$  be a plausible equilibrium of the new biased game. Observe first that  $s_2' \in BR(\psi_2'(s_1')) = BR(s_1^*)$ . This implies that  $s_2' > s_2^*$ , and, thus, due to the monotonicity of  $\psi_1^*$ , we have:  $\psi_1^*(s_2') \geq \psi_1^*(s_2^*)$ . We consider two cases:

1. If  $\psi_1^*(s_2') > \psi_1^*(s_2^*)$ , then the strategic complementarity of player 1's payoff implies that  $s_1' \geq \min(BR(\psi_1^*(s_2'))) \geq \max(BR(\psi_1^*(s_2^*))) \geq s_1^*$ , and, this, in turn, implies that player 2 strictly gains from his deviation:  $\pi_2(s_1', s_2') \geq \pi_2(s_1', s_2^*) > \pi_2(s_1^*, s_2^*)$ , a contradiction.
2. If  $\psi_1^*(s_2') = \psi_1^*(s_2^*)$ , then  $(s_1^*, s_2') \in PNE(G_{(\psi_1^*, \psi_2')})$  and  $\pi_2(s_1^*, s_2') > \pi_2(s_1^*, s_2^*)$ , which contradicts that  $((\psi_1^*, \psi_2^*), (s_1^*, s_2^*))$  is a BBE.

Next we show underinvestment of player 1. Assume to the contrary that  $s_1^* > \max(BR(s_2^*))$ . Consider a deviation of player 1 to a blind belief that the opponent always plays strategy  $s_2^*$  (i.e.,  $\psi_1' \equiv s_2^*$ ). Let  $(s_1', s_2') \in PNE(G_{(\psi_1', \psi_2^*)})$  be a plausible equilibrium of the new biased game. Observe first that  $s_1' \in BR(\psi_1'(s_2')) = BR(s_2^*)$ . This implies that  $s_1' < s_1^*$  and, thus, due to the monotonicity of  $\psi_2^*$ , we have:  $\psi_2^*(s_1') \leq \psi_2^*(s_1^*)$ . We consider two cases:

1. If  $\psi_2^*(s_1') < \psi_2^*(s_1^*)$ , then the strategic substitutability of player 2's payoff implies that  $s_2' \geq \min(BR(\psi_2^*(s_1'))) \geq \max(BR(\psi_2^*(s_1^*))) \geq s_2^*$ , and this, in turn, implies that player 1 strictly gains from his deviation:  $\pi_1(s_1', s_2') \geq \pi_1(s_1', s_2^*) > \pi_1(s_1^*, s_2^*)$ , a contradiction.
2. If  $\psi_2^*(s_1') = \psi_2^*(s_1^*)$ , then  $(s_1', s_2^*) \in PNE(G_{(\psi_1', \psi_2^*)})$  and  $\pi_1(s_1', s_2^*) > \pi_1(s_1^*, s_2^*)$ , which contradicts that  $((\psi_1^*, \psi_2^*), (s_1^*, s_2^*))$  is a BBE.

**Part 2:** Assume that strategy profile  $(s_1^*, s_2^*)$  satisfies I, II, and III. Let  $s_1^e = \min(BR^{-1}(s_2^*))$ . For each strategy  $s_2 < s_2^*$  define  $X(s_2)$  as the set of strategies  $s_2'$  for which player 2 is worse off (relative to  $\pi_2(s_1^*, s_2^*)$ ) if he plays strategy  $s_2$ , while player 1 plays a best-reply to  $s_2'$ . Formally:

$$X_{s^*}(s_2) = \{s_2' \in S_2 | \pi_2(s_1, s_2) \leq \pi_2(s_1^*, s_2') \quad \forall s_1 \in BR(s_2')\}.$$

The assumption that  $\pi_2(s_1^*, s_2^*) > \tilde{M}_2^U$  implies that  $X_{s^*}(s_2)$  is nonempty for each  $s_2 \in S_2$ . The assumption of strategic complements of player 1's payoff implies that  $X_{s^*}(s_2)$  is an interval starting at  $\min(S_2)$ . Let  $\phi_{s^*}(s_2) = \sup(X_{s^*}(s_2))$ . The assumption that the payoff function is continuously twice differentiable implies that  $\phi_{s^*}(s_2)$  is continuous. The assumption that  $s_1^e = \min(BR^{-1}(s_2^*))$  implies that  $\lim_{s_2 \nearrow s_2^*} (\phi_{s^*}(s_2)) = s_1^e$ . These observations imply that there exists a monotone biased belief  $\psi_1^*$  satisfying (1)  $\psi_1^*(s_2) = s_1^e$  and (2)  $\psi_1^*(s_2) \leq \phi_{s^*}(s_2)$  for each  $s_2 < s_2^*$  with an equality only if  $\phi_{s^*}(s_2) = \min(S_2)$ .

Let  $s_2^e = \max(BR^{-1}(s_1^*))$ . For each strategy  $s_1 > s_1^*$  define  $X(s_1)$  as the set of strategies  $s_1' \in S_1$  for which player 1 is worse off (relative to  $\pi_2(s_1^*, s_2^*)$ ) if he plays strategy  $s_1$ , while player 2 plays a best-reply to  $s_1'$ . Formally:

$$X_{s^*}(s_1) = \{s_1' \in S_1 | \pi_1(s_1, s_2) \leq \pi_1(s_1', s_2^*) \quad \forall s_2 \in BR(s_1')\}.$$

The assumption that  $\pi_1(s_1^*, s_2^*) > \tilde{M}_1^U$  implies that  $X_{s^*}(s_1)$  is nonempty for each  $s_1 \in S_1$ . The assumption of strategic substitutes of player 2's payoff implies that  $X_{s^*}(s_1)$  is an interval ending

at  $\max(S_1)$ . Let  $\phi_{s^*}(s_1) = \inf(X_{s^*}(s_1))$ . The assumption that the payoff function is continuously twice differentiable implies that  $\phi_{s^*}(s_1)$  is continuous. The assumption that  $s_2^e = \max(BR^{-1}(s_1^*))$  implies that  $\lim_{s_1 \searrow s_1^*} (\phi_{s^*}(s_1)) = s_1^e$ . These observations imply that there exists a monotone biased belief  $\psi_2^*$  satisfying (1)  $\psi_2^*(s_1) = s_1^e$  and (2)  $\psi_2^*(s_1) \geq \phi_{s^*}(s_1)$  for each  $s_1 > s_1^*$  with an equality only if  $\phi_{s^*}(s_1) = \max(S_1)$ .

We now show that these properties of  $(\psi_1^*, \psi_2^*)$  imply that  $((\psi_1^*, \psi_2^*), (s_1^*, s_2^*))$  is a BBE. Consider a deviation of player 2 into an arbitrary biased belief  $\psi_2'$ . For each  $s_2' \geq s_2^*$ , and each  $(s_1', s_2') \in PNE(G_{(\psi_1^*, \psi_2')})$ , the fact that  $\psi_1^*(s_2') = \psi_1^*(s_2^*)$  implies that  $s_1' = s_1^*$ , and due to assumption (III) of the overinvestment of player 2 and the concavity of the payoff function, we have  $\pi_2(s_1', s_2') = \pi_2(s_1', s_2^*) \leq \pi_2(s_1^*, s_2^*)$ . For each  $s_2' < s_2^*$ , and each  $(s_1', s_2') \in NE(G_{(\psi_1^*, \psi_2')})$ , the fact that  $\psi_1^*(s_2') \leq \phi_{s^*}(s_2')$  with an equality only if  $\phi_{s^*}(s_2') = \min(S_2)$  implies that  $\pi_2(s_1', s_2') \leq \pi_2(s_1^*, s_2^*)$ . This shows that player 2 cannot gain from his deviation.

Finally, consider a deviation of player 1 to an arbitrary biased belief  $\psi_1'$ . For each  $s_1' \leq s_1^*$ , and each  $(s_1', s_2') \in PNE(G_{(\psi_1', \psi_2^*)})$ , the fact that  $\psi_2^*(s_1') = \psi_2^*(s_1^*)$  implies that  $s_2' = s_2^*$ , and due to assumption (III) of the underinvestment of player 1 and the concavity of the payoff function, we have  $\pi_1(s_1', s_2') = \pi_1(s_1', s_2^*) \leq \pi_1(s_1^*, s_2^*)$ . For each  $s_1' > s_1^*$ , and each  $(s_1', s_2') \in NE(G_{(\psi_1', \psi_2^*)})$ , the fact that  $\psi_2^*(s_1') \geq \phi_{s^*}(s_1')$  with an equality only if  $\phi_{s^*}(s_1') = \max(S_1)$  implies that  $\pi_1(s_1', s_2') \leq \pi_1(s_1^*, s_2^*)$ . This shows that player 1 cannot gain from his deviation, which implies that  $((\psi_1^*, \psi_2^*), (s_1^*, s_2^*))$  is a BBE.

## F.9 Proof of Corollary 4 (Pessimism in Games with Strategic Opposites)

The proof is analogous to the proof of Corollary 3, and is presented for completeness.

Assume to the contrary that  $\psi_i^*(s_j^*) > s_j^*$  for some player  $i$ . Assume first that  $\psi_2^*(s_1^*) > s_1^*$ ; then Lemma 4 implies that  $\max(BR(\psi_2^*(s_1^*))) \leq \min(BR(s_1^*))$  with an equality only if

$$\max(BR(\psi_2^*(s_1^*))) \in \{\min(S_2), \max(S_2)\}.$$

Part 1 of Proposition 6 and the definition of a monotone BBE imply that

$$\max(BR(\psi_2^*(s_1^*))) \geq s_2^* \geq \min(BR(s_1^*)).$$

The previous inequalities jointly imply that

$$\max(BR(\psi_2^*(s_1^*))) = s_2^* = \min(BR(s_1^*)) \in \{\min(S_2), \max(S_2)\},$$

which contradicts the assumption that  $s_2^* \notin \{\min(S_2), \max(S_2)\}$ .

We are left with the case of  $\psi_1^*(s_2^*) > s_2^*$ ; then Lemma 2 implies that  $\min(BR(\psi_1^*(s_2^*))) \geq \max(BR(s_2^*))$  with an equality only if

$$\min(BR(\psi_1^*(s_2^*))) \in \{\min(S_1), \max(S_1)\}.$$

Part 1 of Proposition 6 and the definition of a monotone BBE imply that

$$\min(BR(\psi_1^*(s_2^*))) \leq s_1^* \leq \max(BR(s_2^*)).$$

The previous inequalities jointly imply that

$$\min(BR(\psi_1^*(s_2^*))) = s_1^* = \max(BR(s_2^*)) \in \{\min(S_1), \max(S_1)\},$$

which contradicts the assumption that  $s_1^* \notin \{\min(S_1), \max(S_1)\}$ .